

AD-A160 136

COMPARING THE EFFECTIVENESS OF SOFTWARE TESTING  
STRATEGIES(U) MARYLAND UNIV COLLEGE PARK DEPT OF  
COMPUTER SCIENCE V R BASILI ET AL. MAY 85 CS-TR-1501  
AFOSR-TR-85-0005 F49620-80-C-0001

1/1

UNCLASSIFIED

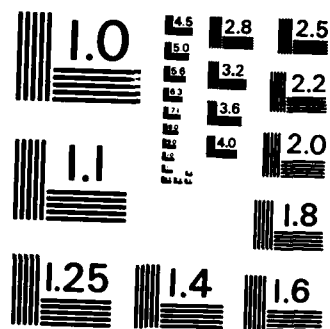
F/G 9/2

NL

END

FILED

DTIC



MICROCOPY RESOLUTION TEST CHART  
NATIONAL BUREAU OF STANDARDS-1963-A

AFOSR-TR- 85-0805

2

AD-A160 136

Technical Report TR-1501

May 1985

Comparing the Effectiveness of  
Software Testing Strategies

Victor R. Basili  
Richard W. Selby, Jr.

Department of Computer Science  
University of Maryland  
College Park

COMPUTER SCIENCE  
TECHNICAL REPORT SERIES



DTIC  
ELECTE  
OCT 11 1985

S

D

B

DTIC FILE COPY

UNIVERSITY OF MARYLAND  
COLLEGE PARK, MARYLAND

20742

Approved for public release;  
distribution unlimited.

85 10 11 079

Technical Report TR-1501

May 1985

Comparing the Effectiveness of  
Software Testing Strategies

Victor R. Basili  
Richard W. Selby, Jr.

Department of Computer Science  
University of Maryland  
College Park

DTIC  
ELECTE  
OCT 11 1985  
S B D

KEYWORDS:

software testing, functional testing, structural testing, code reading,  
off-line software review, empirical study, methodology evaluation,  
software measurement

Research supported in part by the Air Force Office of Scientific Research  
Contract AFOSR-F49620-80-C-001 and the National Aeronautics and Space  
Administration Grant NSG-5123 to the University of Maryland. Computer  
support provided in part by the facilities of NASA/Goddard Space Flight  
Center and the Computer Science Center at the University of Maryland.

**DISTRIBUTION STATEMENT A**

Approved for public release;  
Distribution Unlimited

AIR FORCE OFFICE OF SCIENTIFIC RESEARCH (AFOSR)  
NOTICE OF RESEARCH RESULTS

TR-1501  
May 1985  
DTIC  
MATCHED  
Chief, Technical Information Division

## ABSTRACT

This study compares the strategies of code reading, functional testing, and structural testing in three aspects of software testing: fault detection effectiveness, fault detection cost, and classes of faults detected. Thirty two professional programmers and 42 advanced students applied the three techniques to four unit-sized programs in a fractional factorial experimental design. The major results of this study are the following.

- 1) With the professional programmers, code reading detected more software faults and had a higher fault detection rate than did functional or structural testing, while functional testing detected more faults than did structural testing, but functional and structural testing were not different in fault detection rate.
- 2) In one advanced student subject group, code reading and functional testing were not different in faults found, but were both superior to structural testing, while in the other advanced student subject group there was no difference among the techniques.
- 3) With the advanced student subjects, the three techniques were not different in fault detection rate.
- 4) Number of faults observed, fault detection rate, and total effort in detection depended on the type of software tested.
- 5) Code reading detected more interface faults than did the other methods.
- 6) Functional testing detected more control faults than did the other methods.
- 7) When asked to estimate the percentage of faults detected, code readers gave the most accurate estimates while functional testers gave the least accurate estimates.



Distribution For		<input checked="" type="checkbox"/>
Available for		
Distribution		
Availability		
Dist		
A-1		

## Table of Contents

1 Introduction .....	1
2 Testing Techniques .....	1
2.1 Investigation Goals .....	2
3 Empirical Study .....	2
3.1 Iterative Experimentation .....	3
3.2 Subject and Program/Fault Selection .....	3
3.2.1 Subjects .....	3
3.2.2 Programs .....	5
3.2.3 Faults .....	6
3.2.3.1 Fault Origin .....	6
3.2.3.2 Fault Classification .....	7
3.2.3.3 Fault Description .....	8
3.3 Experimental Design .....	8
3.3.1 Independent and Dependent Variables .....	8
3.3.2 Analysis of Variance Model .....	9
3.4 Experimental Operation .....	12
4 Data Analysis .....	14
4.1 Fault Detection Effectiveness .....	14
4.1.1 Data Distributions .....	14
4.1.2 Number of Faults Detected .....	15
4.1.3 Percentage of Faults Detected .....	15
4.1.4 Dependence on Software Type .....	16
4.1.5 Observable vs. Observed Faults .....	17
4.1.6 Dependence on Program Coverage .....	17
4.1.7 Dependence on Programmer Expertise .....	18
4.1.8 Accuracy of Self-Estimates .....	18
4.1.9 Dependence on Interactions .....	18
4.1.10 Summary of Fault Detection Effectiveness .....	19
4.2 Fault Detection Cost .....	19
4.2.1 Data Distributions .....	20
4.2.2 Fault Detection Rate and Total Time .....	20
4.2.3 Dependence on Software Type .....	21
4.2.4 Computer Costs .....	22
4.2.5 Dependence on Programmer Expertise .....	22

4.2.6 Dependence on Interactions .....	23
4.2.7 Relationships Between Fault Detection Effectiveness and Cost .....	23
4.2.8 Summary of Fault Detection Cost .....	23
4.3 Characterization of Faults Detected .....	24
4.3.1 Omission vs. Commission Classification .....	24
4.3.2 Six-Part Fault Classification .....	25
4.3.3 Observable Fault Classification .....	25
4.3.4 Summary of Characterization of Faults Detected .....	26
5 Conclusions .....	26
6 Acknowledgement .....	28
7 Appendices .....	29
7.1 Appendix A. The Specifications for the Programs .....	29
7.2 Appendix B. The Source Code for the Programs .....	32
8 References .....	50

## 1. Introduction

The processes of software testing and defect detection continue to challenge the software community. Even though the software testing and defect detection activities are inexact and inadequately understood, they are crucial to the success of a software project. The controlled study presented addresses the uncertainty of how to test software effectively. In this investigation, common testing techniques were applied to different types of software by subjects that had a wide range of professional experience. This work is intended to characterize how testing effectiveness relates to several factors: testing technique, software type, fault type, tester experience, and any interactions among these factors. This examination extends previous work by incorporating different testing techniques and a greater number of persons and programs, while broadening the scope of issues examined and adding statistical significance to the conclusions.

The following sections describe the testing techniques examined, the investigation goals, the experimental design, operation, analysis, and conclusions.

## 2. Testing Techniques

To demonstrate that a particular program actually meets its specifications, professional software developers currently utilize many different testing methods. Before presenting the goals for the empirical study comparing the popular techniques of code reading, functional testing, and structural testing, a description will be given of the testing strategies and their different capabilities (see Figure 1.). In functional testing, which is a "black box" approach [Howden 80], a programmer constructs test data from the program's specification through methods such as equivalence partitioning and boundary value analysis [Myers 79]. The programmer then executes the program and contrasts its actual behavior with that indicated in the specification. In structural testing, which is a "white box" approach [Howden 78, Howden 81], a programmer inspects the source code and then devises and executes test cases based on the percentage of the program's statements or expressions executed (the "test set coverage") [Stuckl 77]. The structural coverage criteria used was 100% statement coverage. In code reading by stepwise abstraction, a person identifies prime subprograms in the software, determines their functions, and composes these functions to determine a function for the entire program [Mills 72.



Linger, Mills & Witt 79]. The code reader then compares this derived function and the specifications (the intended function). In order to contrast these various strategies, an empirical study has been conducted using the techniques of code reading, functional testing, and structural testing.

## **2.1. Investigation Goals**

The goals of this study comprise three different aspects of software testing: fault detection effectiveness, fault detection cost, and classes of faults detected. An application of the goal/question/metric paradigm [Basili & Selby 84, Basili & Weiss 84] leads to the framework of goals and questions for this study appearing in Figure 2.

The first goal area is performance oriented and includes a natural first question (I.A): which of the techniques detects the most faults in the programs? The comparison between the techniques is being made across programs, each with a different number of faults. An alternate interpretation would then be to compare the percentage of faults found in the programs (question I.A.1). The number of faults that a technique exposes should also be compared; that is, faults that are made observable but not necessarily observed and reported by a tester (I.A.2). Because of the differences in types of software and in testers' abilities, it is relevant to determine whether the number of faults detected is either program or programmer dependent (I.B, I.C). Since one technique may find a few more faults than another, it becomes useful to know how much effort that technique requires (II.A). Awareness of what types of software require more effort to test (II.B) and what types of programmer backgrounds require less effort in fault uncovering (II.C) is also quite useful. If one is interested in detecting certain classes of faults, such as in error-based testing [Foster 80, Valdes & Goel 83], it is appropriate to apply a technique sensitive to that particular type (III.A). Classifying the types of faults that are observable yet go unreported could help focus and increase testing effectiveness (III.B).

## **3. Empirical Study**

Admittedly, the goals stated here are quite ambitious. In no way is it implied that this study can definitively answer all of these questions for all environments. It is intended, however, that the statistically significant analysis presented lends insights into their answers and into the merit and appropriateness of each of the techniques. Note

that this study compares the individual application of the three testing techniques in order to identify their distinct advantages and disadvantages. This approach is a first step toward proposing a composite testing strategy, which possibly incorporates several testing methods. The following sections describe the empirical study undertaken to pursue these goals and questions, including the selection of subjects, programs, and experimental design, and the overall operation of the study.

### **3.1. Iterative Experimentation**

The empirical study consisted of three phases. The first and second phases of the study took place at the University of Maryland in the Falls of 1982 and 1983 respectively. The third phase took place at Computer Sciences Corporation (CSC - Silver Spring, MD) and NASA Goddard Space Flight Center (Greenbelt, MD) in the Fall of 1984. The sequential experimentation supported the iterative nature of the learning process, and enabled the initial set of goals and questions to be expanded and resolved by further analysis. The goals were further refined by discussions of the preliminary results [Selby 83, Selby 84]. These three phases enabled the pursuit of result reproducibility across environments having subjects with a wide range of experience.

### **3.2. Subject and Program/Fault Selection**

A primary consideration in this study was to use a realistic testing environment to assess the effectiveness of these different testing strategies, as opposed to creating a best possible testing situation [Hetzl 76]. Thus, 1) the subjects for the study were chosen to be representative of different levels of expertise, 2) the programs tested correspond to different types of software and reflect common programming style, and 3) the faults in the programs were representative of those frequently occurring in software. Sampling the subjects, programs, and faults in this manner is intended to evaluate the testing methods reasonably, and to facilitate the generalization of the results to other environments.

#### **3.2.1. Subjects**

The three phases of the study incorporated a total of 74 subjects: the individual phases had 29, 13, and 32 subjects respectively. The subjects were selected, based on

several criteria, to be representative of three different levels of computer science expertise: advanced, intermediate, and junior. The number of subjects in each level of expertise for the different phases appears in Figure 3.

The 42 subjects in the first two phases of the study were the members of the upper level "Software Design and Development" course at the University of Maryland in the Falls of 1982 and 1983. The individuals were either upper-level computer science majors or graduate students; some were working part-time and all were in good academic standing. The topics of the course included structured programming practices, functional correctness, top-down design, modular specification and design, step-wise refinement, and PDL, in addition to the presentation of the techniques of code reading, functional testing, and structural testing. The references for the testing methods were [Mills 75, Fagan 76, Myers 79, Howden 80], and the lectures were presented by V. R. Basili and F. T. Baker. The subjects from the University of Maryland spanned the intermediate and junior levels of computer science expertise. The assignment of individuals to levels of expertise was based on professional experience and prior academic performance in relevant computer science courses. The individuals in the first and second phases had overall averages of 1.7 (SD = 1.7) and 1.5 (SD = 1.5) years of professional experience. The nine intermediate subjects in the first phase had from 2.8 to 7 years of professional experience (average of 3.9 years, SD = 1.3), and the four in the second phase had from 2.3 to 5.5 years of professional experience (average of 3.2, SD = 1.5). The twenty junior subjects in the first phases and the nine in the second phase both had from 0 to 2 years professional experience (averages of 0.7, SD = 0.6, and 0.8, SD = 0.8, respectively).

The 32 subjects in the third phase of the study were programming professionals from NASA and Computer Sciences Corporation. These individuals were mathematicians, physicists, and engineers that develop ground support software for satellites. They were familiar with all three testing techniques, but had used functional testing primarily. A four hour tutorial on the testing techniques was conducted for the subjects by R. W. Selby. This group of subjects, examined in the third phase of the experiment, spanned all three expertise levels and had an overall average of 10.0 (SD = 5.7) years professional experience. Several criteria were considered in the assignment of subjects to

expertise levels, including years of professional experience, degree background, and their manager's suggested assignment. The eight advanced subjects ranged from 9.5 to 20.5 years professional experience (average of 15.0, SD = 4.1). The eleven intermediate subjects ranged from 3.5 to 17.5 years experience (average of 10.9, SD = 4.9). The thirteen junior subjects ranged from 1.5 to 13.5 years experience (average of 6.1, SD = 4.4).

### 3.2.2. Programs

The experimental design enables the distinction of the testing techniques while allowing for the effects of the different programs being tested. The four programs used in the investigation were chosen to be representative of several different types of software. The programs were selected specially for the study and were provided to the subjects for testing; the subjects did not test programs that they had written. All programs were written in a high-level language with which the subjects were familiar. The three programs tested in the CSC/NASA phase were written in FORTRAN, and the programs tested in the University of Maryland phases were written in the Simpl-T structured programming language [Basill & Turner 78].<sup>1</sup> The four programs tested were  $P_1$ ) a text processor,  $P_2$ ) a mathematical plotting routine,  $P_3$ ) a numeric abstract data type, and  $P_4$ ) a database maintainer. The programs are summarized in Figure 4. There exists some differentiation in size, and the programs are a realistic size for unit testing. Each of the subjects tested three programs, but a total of four programs was used across the three phases of the study. The programs tested in each of the three phases of the study appear in Figure 5. The specifications for the programs appear in Appendix A, and their source code appears in Appendix B.

The first program is a text formatting program, which also appeared in [Myers 78]. A version of this program, originally written by [Naur 69] using techniques of program correctness proofs, was analyzed in [Goodenough & Gerhart 75]. The second program is a mathematical plotting routine. This program was written by R. W. Selby, based roughly on a sample program in [Jensen & Wirth 74]. The third program is a numeric

---

<sup>1</sup> Simpl-T is a structured language that supports several string and file handling primitives, in addition to the usual control flow constructs available, for example, in Pascal.

data abstraction consisting of a set of list processing utilities. This program was submitted for a class project by a member of an intermediate level programming course at the University of Maryland. [McMullin & Gannon 80]. The fourth program is a maintainer for a database of bibliographic references. This program was analyzed in [Hetzel 76], and was written by a systems programmer at the University of North Carolina computation center.

Note that the source code for the programs contains no comments. This creates a worst-case situation for the code readers. In an environment where code contained helpful comments, performance of code readers would likely improve, especially if the source code contained as comments the intermediate functions of the program segments. In an environment where the comments were at all suspect, they could then be ignored.

### 3.2.3. Faults

The faults contained in the programs tested represent a reasonable distribution of faults that commonly occur in software [Weiss & Basili 85, Basili & Perricone 84]. All the faults in the database maintainer and the numeric abstract data type were made during the actual development of the programs. The other two programs contain a mix of faults made by the original programmer and faults seeded in the code. The programs contained a total of 34 faults; the text formatter had nine, the plotting routine had six, the abstract data type had seven, and the database maintainer had twelve.

#### 3.2.3.1. Fault Origin

The faults in the text formatter were preserved from the article in which it appeared [Myers 78], except for some of the more controversial ones [Calliau & Rubin 79]. In the mathematical plotter, faults made during program translation were supplemented by additional representative faults. The faults in the abstract data type were the original ones made by the program's author during the development of the program. The faults in the database maintainer were recorded during the development of the program, and then reinserted into the program. The next section describes a classification of the different types of faults in the programs. Note that this investigation of the fault detecting ability of these techniques involves only those types occurring in the source code, not other types such as those in the requirements or the specifications.

### 3.2.3.2. Fault Classification

The faults in the programs are classified according to two different abstract classification schemes [Basili & Perricone 84]. One fault categorization method separates faults of omission from faults of commission. Faults of commission are those faults present as a result of an incorrect segment of existing code. For example, the wrong arithmetic operator is used for a computation in the right-hand-side of an assignment statement. Faults of omission are those faults present as a result of a programmer's forgetting to include some entity in a module. For example, a statement is missing from the code that would assign the proper value to a variable.

A second fault categorization scheme partitions software faults into the six classes of 1) initialization, 2) computation, 3) control, 4) interface, 5) data, and 6) cosmetic. Improperly initializing a data structure constitutes an initialization fault. For example, assigning a variable the wrong value on entry to a module. Computation faults are those that cause a calculation to evaluate the value for a variable incorrectly. The above example of a wrong arithmetic operator in the right-hand-side of an assignment statement would be a computation fault. A control fault causes the wrong control flow path in a program to be taken for some input. An incorrect predicate in an IF-THEN-ELSE statement would be a control fault. Interface faults result when a module uses and makes assumptions about entities outside the module's local environment. Interface faults would be, for example, passing an incorrect argument to a procedure, or assuming in a module that an array passed as an argument was filled with blanks by the passing routine. A data fault are those that result from the incorrect use of a data structure. For example, incorrectly determining the index for the last element in an array. Finally, cosmetic faults are clerical mistakes when entering the program. A spelling mistake in an error message would be a cosmetic fault.

Interpreting and classifying faults in software is a difficult and inexact task. The categorization process often requires trying to recreate the original programmer's misunderstanding of the problem [Johnson, Draper & Soloway 83]. The above two fault classification schemes attempt to distinguish among different reasons that programmers make faults in software development. They were applied to the faults in the programs

In a consistent interpretation; it is certainly possible that another analyst could have interpreted them differently. The separate application of each of the two classification schemes to the faults categorized them in a mutually exclusive and exhaustive manner. Figure 6 displays the distribution of faults in the programs according to these schemes.

### 3.2.3.3. Fault Description

The faults in the programs are described in Figure 7. There have been various efforts to determine a precise counting scheme for "defects" in software [Gloss-Soler 79, IEEE 83]. According to the explanations given, a software "fault" is a specific manifestation in the source code of a programmer "error." For example, due to a misconception or document discrepancy, a programmer commits an "error" (in his/her head) that may result in more than one "fault" in a program. Using this interpretation, software "faults" reflect the correctness, or lack thereof, in a program. The entities examined in this analysis are software faults.

### 3.3. Experimental Design

The experimental design applied for each of the three phases of the study was a fractional factorial design [Cochran & Cox 50, Box, Hunter, & Hunter 78]. This experimental design distinguishes among the testing techniques, while allowing for variation in the ability of the particular individual testing or in the program being tested. Figure 8 displays the fractional factorial design appropriate for the third phase of the study. Subject  $S_1$  is in the advanced expertise level, and he structurally tested program  $P_1$ , functionally tested program  $P_3$ , and code read program  $P_4$ . Notice that all of the subjects tested each of the three programs and used each of the three techniques. Of course, no one tests a given program more than once. The design appropriate for the third phase is discussed in the following paragraphs, with the minor differences between this design and the ones applied in the first two phases being discussed at the end of the section.

#### 3.3.1. Independent and Dependent Variables

The experimental design has the three independent variables of testing technique, software type, and level of expertise. For the design appearing in Figure 8, appropriat

for the third phase of the study, the three main effects have the following levels:

- 1) testing technique: code reading, functional testing, and structural testing
- 2) software type: ( $P_1$ ) text processing, ( $P_2$ ) numeric abstract data type, and ( $P_3$ ) database maintainer
- 3) level of expertise: advanced, intermediate, and junior

Every combination of these levels occurs in the design. That is, programmers in all three levels of expertise applied all three testing techniques on all programs. In addition to these three main effects, a factorial analysis of variance (ANOVA) model supports the analysis of interactions among each of these main effects. Thus, the interaction effects of testing technique \* software type, testing technique \* expertise level, software type \* expertise level, and the three-way interaction of testing technique \* software type \* expertise level are included in the model. There are several dependent variables examined in the study, including number of faults detected, percentage of faults detected, total fault detection time, and fault detection rate. Observations from the on-line methods of functional and structural testing also had as dependent variables number of computer runs, amount of cpu-time consumed, maximum statement coverage achieved, connect time used, number of faults that were observable from the test data, percentage of faults that were observable from the test data, and percentage of faults observable in the from the test data that were actually observed by the tester.

### 3.3.2. Analysis of Variance Model

The three main effects and all the two-way and three-way interactions effects are called fixed effects in this factorial analysis of variance model. The levels of these effects given above represent all levels of interest in the investigation. For example, the effect of testing technique has as particular levels code reading, functional testing, and structural testing; these particular testing techniques are the only ones under comparison in this study. The effect of the particular subjects that participated in this study requires a little different interpretation. The subjects examined in the study were random samples of programmers from the large population of programmers at each of the levels of expertise. Thus, the effect of the subjects on the various dependent variables is a random variable, and this effect therefore is called a random effect. If the samples exam-



ined are truly representative of the population of subjects at each expertise level, the inferences from the analysis can then be generalized across the whole population of subjects at each expertise level, not just across the particular subjects in the sample chosen. Since this analysis of variance model contains both fixed and random effects, it is called a mixed model. The actual ANOVA model for the design appearing in Figure 8 is given below.

$$T_{ijkl} = \mu + \alpha_i + \beta_j + \gamma_k + \delta_{kl} + \alpha\beta_{ij} + \alpha\gamma_{ik} + \beta\gamma_{jk} + \alpha\beta\gamma_{ijk} + \epsilon_{ijkl}$$

where

$T_{ijkl}$  is the observed response from subject  $l$  of experience level  $k$  using testing technique  $i$  on program  $j$

$\mu$  is the overall mean response

$\alpha_i$  is the main effect of testing technique  $i$  ( $i = 1, 2, 3$ )

$\beta_j$  is the main effect of program  $j$  ( $j = 1, 3, 4$ )

$\gamma_k$  is the main effect of expertise level  $k$  ( $k = 1, 2, 3$ )

$\delta_{kl}$  is the random effect of subject  $l$  within expertise level  $k$ , a random variable ( $l = 1, 2, \dots, 32; k = 1, 2, 3$ )

$\alpha\beta_{ij}$  is the interaction effect of testing technique  $i$  with program  $j$  ( $i = 1, 2, 3; j = 1, 3, 4$ )

$\alpha\gamma_{ik}$  is the interaction effect of testing technique  $i$  with expertise level  $k$  ( $i = 1, 2, 3; k = 1, 2, 3$ )

$\beta\gamma_{jk}$  is the interaction effect of program  $j$  with expertise level  $k$  ( $j = 1, 3, 4; k = 1, 2, 3$ )

$\alpha\beta\gamma_{ijk}$  is the interaction effect of testing technique  $i$  with program  $j$  with experience level  $k$  ( $i = 1, 2, 3; j = 1, 3, 4; k = 1, 2, 3$ )

$\epsilon_{ijkl}$  is the experimental error for each observation, a random variable

The F tests of hypotheses on all the fixed effects mentioned above use the error (residual) mean square in the denominator, except for the test of the expertise level effect. The expected mean square for the expertise level effect contains a component for the actual variance of subjects within expertise level. In order to select the appropriate error term for the denominator of the expertise level F test, the mean square for the effect of subjects nested within expertise level is chosen. The parameters for the random effect of subjects within expertise level are assumed to be drawn from a normally distributed random process with mean zero and common variance. The experimental error terms are assumed to have mean zero and common variance.

The fractional factorial design applied in the first two phases of the analysis differed slightly from the one presented above for the third phase.<sup>2</sup> In the third phase of the study, programs  $P_1$ ,  $P_3$ , and  $P_4$  were tested by subjects in three levels of expertise. In both phases one and two, there were only subjects from the levels of intermediate and junior expertise. In phase one, programs  $P_1$ ,  $P_3$ , and  $P_2$  were tested. In phase two, the programs tested were  $P_1$ ,  $P_2$ , and  $P_4$ . The only modifications necessary to the above explanation for phases one and two are 1) eliminating the advanced expertise level, 2) changing the program P subscripts appropriately, and 3) leaving out the three way interaction term in phase two, because of the reduced number of subjects. In all three of the phases, all subjects used each of the three techniques and tested each of the three programs for that phase. Also, within all three phases, all possible combinations of expertise level, testing techniques, and programs occurred.

The order of presentation of the testing techniques was randomized among the subjects in each level of expertise in each phase of the study. However, the integrity of the results would have suffered if each of the programs in a given phase was tested at different times by different subjects. Note that each of the testing sessions took place on a different day because of the amount of effort required. If different programs would have been tested on different days, any discussion about the programs among subjects

---

<sup>2</sup> Although the data from all the phases can be analyzed together, the number of empty cells resulting from not having all three experience levels and all four programs in all phases limits the number of parameters that can be estimated and causes non-unique Type IV partial sums of squares.

between testing sessions would have affected the future performance of others. Therefore, all subjects in a phase tested the same program on the same day. The actual order of program presentation was the order in which the programs are listed in the previous paragraph.

### 3.4. Experimental Operation

Each of the three phases were broken into five distinct pieces: training, three testing sessions, and a follow-up session. All groups of subjects were exposed to a similar amount of training on the testing techniques before the study began. As mentioned earlier, the University of Maryland subjects were enrolled in the "Software Design and Development" course, and the NASA/CSC subjects were given a four-hour tutorial. Background information on the subjects was captured through a questionnaire. Elementary exercises followed by a pretest covering all techniques were administered to all subjects after the training and before the testing sessions. Reasonable effort on the part of the University of Maryland subjects was enforced by their being graded on the work and by their needing to use the techniques in a major class project. Reasonable effort on the part of the NASA/CSC subjects was certain because of their desire for the study's outcome to improve their software testing environment. All subjects groups were judged highly motivated during the study. The subjects were all familiar with the editors, terminals, machines, and the programs' implementation language.

The individuals were requested to use the three testing techniques to the best of their ability. Every subject participated in all three testing sessions of his/her phase, using all techniques but each on a separate program. The individuals using code reading were each given the specification for the program and its source code. They were then asked to apply the methods of code reading by stepwise abstraction to detect discrepancies between the program's abstracted function and the specification. The functional testers were each given a specification and the ability to execute the program. They were asked to perform equivalence partitioning and boundary value analysis to select a set of test data for the program. Then they executed the program on this collection of test data, and inconsistencies between what the program actually performed and what they thought the specification said it should perform were noted. The struc-

tural testers were given the source code for the program, the ability to execute it, and a description of the input format for the program. The structural testers were asked to examine the source and generate a set of test cases that cumulatively execute 100% of the program's statements. When the subjects were applying an on-line technique, they generated and executed their own test data; no test data sets were provided. The programs were invoked through a test driver that supported the use of multiple input data sets. This test driver, unbeknown to the subjects, drained off the input cases submitted to the program for the experimenter's later analysis; the programs could only be accessed through a test driver.

A structural coverage tool calculated the actual statement coverage of the test set and which statements were left unexecuted for the structural testers. After the structural testers generated a collection of test data that met (or almost met) the 100% coverage criteria, no further execution of the program or reference to the source code was allowed. They retained the program's output from the test cases they had generated. These testers were then provided with the program's specification. Now that they knew what the program was intended to do, they were asked to contrast the program's specification with the behavior of the program on the test data they derived. This scenario for the structural testers was necessary so that "observed" faults could be compared.

At the end of each of the testing sessions, the subjects were asked to give a reasonable estimate of the amount of time spent detecting faults with a given testing technique. The University of Maryland subjects were assured that this had nothing to do with the grading of the work. There seemed to be little incentive for the subjects in any of the groups not to be truthful. At the completion of each testing session, the NASA/CSC subjects were also asked what percentage of the faults in the program that they thought were uncovered. After all three testing sessions in a given phase were completed, the subjects were requested to critique and evaluate the three testing techniques regarding their understandability, naturalness, and effectiveness. The University of Maryland subjects submitted a written critique, while a two hour debriefing forum was conducted for the NASA/CSC individuals. In addition to obtaining the impressions of the individuals, these follow-up procedures gave an understanding of how well the

subjects were comprehending and applying the methods. These final sessions also afforded the participants an opportunity to comment on any particular problems they had with the techniques or in applying them to the given programs.

#### 4. Data Analysis

The analysis of the data collected from the various phases of the experiment is presented according to the goal and question framework discussed earlier.

##### 4.1. Fault Detection Effectiveness

The first goal area addresses the fault detection effectiveness of each of the techniques. Figure 9 presents a summary of the measures that were examined to pursue this goal area. A brief description of each measure is as follows - (\*) means only relevant for on-line testing. a) # Faults detected - the number of faults detected by a subject applying a given testing technique on a given program. b) % Faults detected - the percentage of a program's faults that a subject detected by applying a testing technique to the program. c) # Faults observable (\*) - the number of faults that were observable from the program's behavior given the input data submitted. d) % Faults observable (\*) - the percentage of a program's faults that were observable from the program's behavior given the input data submitted. e) % Detected/observable (\*) - the percentage of faults observable from the program's behavior on the given input set that were actually observed by a subject. f) % Faults felt found - a subject's estimate of the percentage of a program's faults that he/she thought were detected by his/her testing. g) Maximum statement coverage (\*) - the maximum percentage of a program's statements that were executed in a set of test cases.

##### 4.1.1. Data Distributions

The actual distribution of the number of faults observed by the subjects appears in Figure 10, broken down by phase. From Figures 9 and 10, the large variation in performance among the subjects is clearly seen. The mean number of faults detected by the subjects is displayed in Figure 11, broken down by technique, program, expertise level, and phase.

#### 4.1.2. Number of Faults Detected

The first question under this goal area asks which of the testing techniques detected the most faults in the programs. The overall F-test of the techniques detecting an equal number of faults in the programs is rejected in the first and third phases of the study ( $\alpha < .024$  and  $\alpha < .0001$ , respectively; not rejected in phase two,  $\alpha > .05$ ). Recall that the phase three data was collected from 32 NASA/CSC subjects, and the phase one data was from 29 University of Maryland subjects. With the phase three data, the contrast of "reading -  $0.5 * (\text{functional} + \text{structural})$ " estimates that the technique of code reading by stepwise abstraction detected 1.24 more faults per program than did either of the other techniques ( $\alpha < .0001$ , c.i. 0.73 - 1.75).<sup>3</sup> Note that code reading performed well even though the professional subjects' primary experience was with functional testing. Also with the phase three data, the contrast of "functional - structural" estimates that the technique of functional testing detected 1.11 more faults per program than did structural testing ( $\alpha < .0007$ , c.i. 0.52 - 1.70). In the phase one data, the contrast of " $0.5 * (\text{reading} + \text{functional}) - \text{structural}$ " estimates that the technique of structural testing detected 1.00 fault less per program than did either reading or functional testing ( $\alpha < .0065$ , c.i. 0.31 - 1.69). In the phase one data, the contrast of "reading - functional" was not statistically different from zero ( $\alpha > .05$ ). The poor performance of structural testing across the phases suggests the inadequacy of using statement coverage criteria. The above pairs of contrasts were chosen because they are linearly independent.

#### 4.1.3. Percentage of Faults Detected

Since the programs tested each had a different number of faults, a question in the earlier goal/question framework asks which technique detected the greatest percentage of faults in the programs. The order of performance of the techniques is the same as above when the percentage of the programs' faults detected are compared. The overall F-tests for phases one and three were rejected as before ( $\alpha < .037$  and  $\alpha < .0001$  respectively; not rejected in phase two,  $\alpha > .05$ ). Applying the same contrasts as above: a) in phase three, reading detected 16.0% more faults per program than did the other tech-

---

<sup>3</sup> The probability of Type I error is reported, the probability of erroneously rejecting the null hypothesis. The abbreviation "c.i." stands for 95% confidence interval.

niques ( $\alpha < .0001$ , c.l. 9.9 - 22.1), and functional detected 11.2% more faults than did structural ( $\alpha < .003$ , c.l. 4.1 - 18.3); b) in phase one, structural detected 13.2% fewer of a program's faults than did the other methods ( $\alpha < .011$ , c.l. 3.5 - 22.9), and reading and functional were not statistically different as before.

#### 4.1.4. Dependence on Software Type

Another question in this goal area queries whether the number or percentage of faults detected depends on the program being tested. The overall F-test that the number of faults detected is not program dependent is rejected only in the phase three data ( $\alpha < .0001$ ). Applying Tukey's multiple comparison on the phase three data reveals that the most faults were detected in the abstract data type, the second most in the text formatter, and the least number of faults were found in the database maintainer (simultaneous  $\alpha < .05$ ). When the percentage of faults found in a program is considered, however, the overall F-tests for the three phases are all rejected ( $\alpha < .027$ ,  $\alpha < .01$ , and  $\alpha < .0001$  in respective order). Tukey's multiple comparison yields the following orderings on the programs (all simultaneous  $\alpha < .05$ ). In the phase one data, the ordering was (data type  $\simeq$  plotter)  $>$  text formatter; that is, a higher percentage of faults were detected in either the abstract data type or the plotter than were found in the text formatter; there was no difference between the abstract data type and the plotter in the percentage found. In the phase two data, the ordering of percentage of faults detected was plotter  $>$  (text formatter  $\simeq$  database maintainer). In the phase three data, the ordering of percentage of faults found in the programs was the same as the number of faults found, abstract data type  $>$  text formatter  $>$  database maintainer. Summarizing the effect of the type of software on the percentage of faults observed: 1) the programs with the highest percentage of their faults detected were the abstract data type and the mathematical plotter, the percentage detected between these two was not statistically different; 2) the programs with the lowest percentage of their faults detected were the text formatter and the database maintainer; the percentage detected between these two was not statistically different in the phase two data, but a higher percentage of faults in the text formatter was detected in the phase three data.

#### 4.1.5. Observable vs. Observed Faults

One evaluation criteria of the success of a software testing session is the number of faults detected. An evaluation criteria of the particular test data generated, however, is the ability of the test data to reveal faults in the program. A test data set's ability to uncover faults in a program can be measured by the number or percentage of a program's faults that are made observable from execution on that input. Distinguishing the faults observable in a program from the faults actually observed by a tester highlights the differences in the activities of test data generation and program behavior examination. As shown in Figure 8, the average number of the programs' faults observable was 68.0% when individuals were either functional testing or structurally testing. Of course, with a nonexecution-based technique such as code reading, 100% of the faults are observable. Test data generated by subjects using the technique of functional testing resulted in 1.4 more observable faults ( $\alpha < .0002$ , c.i. 0.79 - 2.01) than did the use of structural testing in phase one of the study; the percentage difference of functional over structural was estimated at 20.0% ( $\alpha < .0002$ , c.i. 11.2 - 28.8). The techniques did not differ in these two measures in the third phase of the study. However, just considering the faults that were observable from the submitted test data, functional testers detected 18.5% more of these observable faults than did structural testers in the phase three data ( $\alpha < .0016$ , c.i. 8.9 - 28.1); they did not differ in the phase one data. Note that all faults in the programs could be observed in the programs' output given the proper input data. When using the on-line techniques of functional and structural testing, subjects detected 70.3% of the faults observable in the program's output. In order to conduct a successful testing session, faults in a program must be both revealed and subsequently observed.

#### 4.1.6. Dependence on Program Coverage

Another measure of the ability of a test set to reveal a program's faults is the percentage of a program's statements that are executed by the test set. The average maximum statement coverage achieved by the functional and structural testers was 97.0%. The maximum statement coverage from the submitted test data was not statistically different between the functional and structural testers ( $\alpha > .05$ ). Also, there was no correlation between maximum statement coverage achieved and either number or per-



centage of faults found ( $\alpha > .05$ ).

#### 4.1.7. Dependence on Programmer Expertise

A final question in this goal area concerns the contribution of programmer expertise to fault detection effectiveness. In the phase three data from the NASA/CSC professional environment, subjects of advanced expertise detected more faults than did either the subjects of intermediate or junior expertise ( $\alpha < .05$ ). When the percentage of faults detected is compared, however, the advanced subjects performed better than the junior subjects ( $\alpha < .05$ ), but were not statistically different from the intermediate subjects ( $\alpha > .05$ ). The intermediate and junior subjects were not statistically different in any of the three phases of the study in terms of number or percentage faults observed. When several subject background attributes were correlated with the number of faults found, total years of professional experience had a minor relationship (Pearson  $R = .22$ ,  $\alpha < .05$ ). Correspondence of performance with background aspects was examined across all observations, and within each of the phases, including previous academic performance for the University of Maryland subjects. Other than the above, no relationships were found.

#### 4.1.8. Accuracy of Self-Estimates

Recall that the NASA/CSC subjects in the phase three data estimated, at the completion of a testing session, the percentage of a program's faults they thought they had uncovered. This estimation of the number of faults uncovered correlated reasonably well with the actual percentage of faults detected ( $R = .57$ ,  $\alpha < .0001$ ). Investigating further, individuals using the different techniques were able to give better estimates: code readers gave the best estimates ( $R = .79$ ,  $\alpha < .0001$ ), structural testers gave the second best estimates ( $R = .57$ ,  $\alpha < .0007$ ), and functional testers gave the worst estimates (no correlation,  $\alpha > .05$ ). This last observation suggests that the code readers were more certain of the effectiveness they had in revealing faults in the programs.

#### 4.1.9. Dependence on Interactions

There were few significant interactions between the main effects of testing technique, program, and expertise level. In the phase two data, there was an interaction

between testing technique and program in both the number and percentage of faults found ( $\alpha < .0013$ ,  $\alpha < .0014$  respectively). The effectiveness of code reading increased on the text formatter. In the phase three data, there was a slight three-way interaction between testing technique, program, and expertise level for both the number and percentage of faults found ( $\alpha < .05$ ,  $\alpha < .04$  respectively).

#### 4.1.10. Summary of Fault Detection Effectiveness

Summarizing the major results of the comparison of fault detection effectiveness: 1) in the phase three data, code reading detected a greater number and percentage of faults than the other methods, with functional detecting more than structural; 2) in the phase one data, code reading and functional were equally effective, while structural was inferior to both - there were no differences among the three techniques in phase two; 3) the number of faults observed depends on the type of software: the most faults were detected in the abstract data type and the mathematical plotter, the second most in the text formatter, and (in the case of the phase three data) the least were found in the database maintainer; 4) functionally generated test data revealed more observable faults than did structurally generated test data in phase one, but not in phase three; 5) subjects of intermediate and junior expertise were equally effective in detecting faults, while advanced subjects found a greater number of faults than did either group; and 6) self-estimates of faults detected were most accurate from subjects applying code reading, followed by those doing structural testing, with estimates from persons functionally testing having no relationship.

#### 4.2. Fault Detection Cost

The second goal area examines the fault detection cost of each of the techniques. Figure 12 presents a summary of the measures that were examined to investigate this goal area. A brief description of each measure is as follows - (\*) means only relevant for on-line testing. a) # Faults / hour - the number of faults detected by a subject applying a given technique normalized by the effort in hours required, called the fault detection rate. b) Detection time - the total number of hours that a subject spent in testing a program using a technique. c) Cpu-time (\*) - the cpu-time in seconds used during the testing session. d) Normalized cpu-time (\*) - the cpu-time in seconds used during the

testing session, normalized by a factor for machine speed.<sup>4</sup> e) Connect time (\*) - the number of minutes that a individual spent on-line while testing a program. f) # Program runs (\*) - the number of executions of the program test driver; note that the driver supported multiple sets of input data. All of the on-line statistics were monitored by the operating systems of the machines.

#### 4.2.1. Data Distributions

The actual distribution of the fault detection rates for the subjects appears in Figure 13, broken down by phase. Once again, note the many-to-one differential in subject performance. Figure 14 displays the mean fault detection rate for the subjects, broken down by technique, program, expertise level, and phase.

#### 4.2.2. Fault Detection Rate and Total Time

The first question in this goal area asks which testing technique had the highest fault detection rate. The overall F-test of the techniques' having the same fault detection rate was rejected in the phase three data ( $\alpha < .0014$ ), but not in the other two phases ( $\alpha > .05$ ). As before, the two contrasts of "reading -  $0.5 * (\text{functional} + \text{structural})$ " and "functional - structural" were examined to detect differences among the techniques. The technique of code reading was estimated at detecting 1.49 more faults per hour than did the other techniques in the phase three data ( $\alpha < .0003$ , c.i. 0.75 - 2.23). The techniques of functional and structural testing were not statistically different ( $\alpha > .05$ ). Comparing the total time spent in fault detection, the techniques were not statistically different in the phase two and three data; the overall F-test for the phase one data was rejected ( $\alpha < .013$ ). In the phase one data, structural testers spent an estimated 1.08 hours less testing than did the other techniques ( $\alpha < .004$ , c.i. 0.39 - 1.78), while code readers were not statistically different from functional testers. Recall that in phase one, the structural testers observed both a lower number and percentage of the programs' faults than did the other techniques.

---

<sup>4</sup> In the phase three data, testing was done on both a VAX 11/780 and an IBM 4341. As suggested by benchmark comparisons [Church 84], the VAX cpu-times were divided by 1.8 and the IBM cpu-times were divided by 0.9.

#### 4.2.3. Dependence on Software Type

Another question in this area focuses on how fault detection rate depends on software type. The overall F-test that the detection rate is the same for the programs is rejected in the phase one and phase three data ( $\alpha < .01$  and  $\alpha < .0001$  respectively); the detection rate among the programs was not statistically different in phase two. Applying Tukey's multiple comparisons on the phase one data finds that the fault detection rate was greater on the abstract data type than on the plotter, while there was no difference either between the abstract data type and the text formatter or between the text formatter and the plotter (simultaneous  $\alpha < .05$ ). In the phase three data, the fault detection rate was higher in the abstract data type than it was for the text formatter and the database maintainer, with the text formatter and the database maintainer not being statistically different (simultaneous  $\alpha < .05$ ). The overall effort spent in fault detection was different among the programs in phases one and three ( $\alpha < .012$  and  $\alpha < .0001$  respectively), while there was no difference in phase two. In phase one, more effort was spent testing the plotter than the abstract data type, while there was no statistical difference either between the plotter and the text formatter or between the text formatter and the abstract data type (simultaneous  $\alpha < .05$ ). In phase three, more time was spent testing the database maintainer than was spent on either the text formatter or on the abstract data type, with the text formatter not differing from the abstract data type (simultaneous  $\alpha < .05$ ). Summarizing the dependence of fault detection cost on software type, 1) the abstract data type had a higher detection rate and less total detection effort than did either the plotter or the database maintainer, the latter two were not different in either detection rate or total detection time; 2) the text formatter and the plotter did not differ in fault detection rate or total detection effort; 3) the text formatter and the database maintainer did not differ in fault detection rate overall and did not differ in total detection effort in phase two, but the database maintainer had a higher total detection effort in phase three; 4) the text formatter and the abstract data type did not differ in total detection effort overall and did not differ in fault detection rate in phase one, but the abstract data type had a higher detection rate in phase three.

#### 4.2.4. Computer Costs

In addition to the effort spent by individuals in software testing, on-line methods incur machine costs. The machine cost measures of cpu-time, connect time, and the number of runs were compared across the on-line techniques of functional and structural testing in phase three of the study. A nonexecution-based technique such as code reading, of course, incurs no machine time costs. When the machine speeds are normalized (see measure definitions above), the technique of functional testing used 26.0 more seconds of cpu-time than did the technique of structural testing ( $\alpha < .016$ , c.l. 7.0 - 45.0). The estimate of the difference is 29.6 seconds when the cpu-times are not normalized ( $\alpha < .012$ , c.l. 9.0 - 50.2). Individuals using functional testing used 28.4 more minutes of connect time than did those using structural testing ( $\alpha < .004$ , c.l. 11.7 - 45.1). The number of computer runs of a program's test driver was not different between the two techniques ( $\alpha > .05$ ). These results suggest that individuals using functional testing spent more time on-line and used more cpu-time per computer run than did those structurally testing.

#### 4.2.5. Dependence on Programmer Expertise

The relation of programmer expertise to cost of fault detection is another question in this goal section. The expertise level of the subjects had no relation to the fault detection rate in phases two and three ( $\alpha > .05$  for both F-tests). Recall that phase three of the study used 32 professional subjects with all three levels of computer science expertise. In phase one, however, the intermediate subjects detected faults at a faster rate than did the junior subjects ( $\alpha < .005$ ). The total effort spent in fault detection was not different among the expertise levels in any of the phases ( $\alpha > .05$  for all three F-tests). When all 74 subjects are considered, years of professional experience correlates positively with fault detection rate ( $R = .41$ ,  $\alpha < .0002$ ) and correlates slightly negatively with total detection time ( $R = -.25$ ,  $\alpha < .03$ ). These last two observations suggest that persons with more years of professional experience detected the faults faster and spent less total time doing so. Several other subject background measures showed no relationship with fault detection rate or total detection time ( $\alpha < .05$ ). Background measures were examined across all subjects and within the groups of NASA/CSC sub-

jects and University of Maryland subjects.

#### 4.2.6. Dependence on Interactions

There were few significant interactions between the main effects of testing technique, program, and expertise level. There was an interaction between testing technique and software type in terms of fault detection rate and total detection cost for the phase three data ( $\alpha < .003$  and  $\alpha < .007$  respectively). Subjects using code reading on the abstract data type had an increased fault detection rate and a decreased total detection time.

#### 4.2.7. Relationships Between Fault Detection Effectiveness and Cost

There were several correlations between fault detection cost measures and performance measures. Fault detection rate correlated overall with number of faults detected ( $R = .48$ ,  $\alpha < .0001$ ), percentage of faults found ( $R = .48$ ,  $\alpha < .0001$ ), and total detection time ( $R = -.53$ ,  $\alpha < .0001$ ), but not with normalized cpu-time, raw cpu-time, connect time, or number of computer runs ( $\alpha > .05$ ). Total detection time correlated with normalized cpu-time ( $R = .36$ ,  $\alpha < .04$ ) and raw cpu-time ( $R = .37$ ,  $\alpha < .04$ ), but not with connect time, number of runs, number of faults detected, or percentage of faults detected. The number of faults detected in the programs correlated with the amount of machine resources used: normalized cpu-time ( $R = .47$ ,  $\alpha < .007$ ), raw cpu-time ( $R = .52$ ,  $\alpha < .002$ ), and connect time ( $R = .49$ ,  $\alpha < .003$ ), but not with the number of computer runs ( $\alpha > .05$ ). The correlations for percentage of faults detected with machine resources used were similar. Although most of these correlations are minor, they suggest that 1) the higher the fault detection rate, the more faults found and the less time spent in fault detection; 2) fault detection rate had no relationship with use of machine resources; 3) spending more time in detecting faults had no relationship with the amount of faults detected; and 4) the more cpu-time and connect time used, the more faults found.

#### 4.2.8. Summary of Fault Detection Cost

Summarizing the major results of the comparison of fault detection cost: 1) in the phase three data, code reading had a higher fault detection rate than the other methods.

with no difference between functional testing and structural testing; 2) in the phase one and two data, the three techniques were not different in fault detection rate; 3) in the phase two and three data, total detection effort was not different among the techniques, but in phase one less effort was spent for structural testing than for the other techniques, while reading and functional were not different; 4) fault detection rate and total effort in detection depended on the type of software: the abstract data type had the highest detection rate and lowest total detection effort, the plotter and the database maintainer had the lowest detection rate and the highest total detection effort, and the text formatter was somewhere in between depending on the phase; 5) functional testing used more cpu-time and connect time than did structural testing, but they were not different in the number of runs; 6) in phases two and three, subjects across expertise levels were not different in fault detection rate or total detection time, in phase one intermediate subjects had a higher detection rate; and 7) there was a moderate correlation between fault detection rate and years of professional experience across all subjects.

#### 4.3. Characterization of Faults Detected

The third goal area focuses on determining what classes of faults are detected by the different techniques. In the earlier section on the faults in the software, the faults were characterized by two different classification schemes: omission or commission, and initialization, control, data, computation, interface, or cosmetic. The faults detected across all three study phases are broken down by the two fault classification schemes in Figure 15. The entries in the figure are the average percentage (with standard deviations) of faults in a given class observed when a particular technique was being used. Note that when a subject tested a program that had no faults in a given class, he/she was excluded from the calculation of this average.

##### 4.3.1. Omission vs. Commission Classification

When the faults are partitioned according to the omission/commission scheme, there is a distinction among the techniques. Both code readers and functional testers observed more omission faults than did structural testers ( $\alpha < .001$ ), with code readers and functional testers not being different ( $\alpha > .05$ ). Since a fault of omission occurs as a result of some segment of code being left out, you would not expect structurally generat-

ed test data to find such faults. In fact, 44% of the subjects applying structural testing found zero faults of omission when testing a program. A distribution of the faults observed according to this classification scheme appears in Figure 16.

#### 4.3.2. Six-Part Fault Classification

When the faults are divided according to the second fault classification scheme, several differences are apparent. Both code reading and functional testing found more initialization faults than did structural testing ( $\alpha < .05$ ), with code reading and functional testing not being different ( $\alpha > .05$ ). Code reading detected more interface faults than did either of the other methods ( $\alpha < .01$ ), with no difference between functional and structural testing ( $\alpha > .05$ ). This suggests that the code reading process of abstracting and composing program functions across modules must be an effective technique for finding interface faults. Functional testing detected more control faults than did either of the other methods ( $\alpha < .01$ ), with code reading and structural testing not being different ( $\alpha > .05$ ). Recall that the structural test data generation criteria examined is based on determining the execution paths in a program and deriving test data that execute 100% of the program's statements. One would expect that more control path faults would be found by such a technique. However, structural testing did not do as well as functional testing in this fault class. The technique of code reading found more computation faults than did structural testing ( $\alpha < .05$ ), with functional testing not being different from either of the other two methods ( $\alpha > .05$ ). The three techniques were not statistically different in the percentage of faults they detected in either the data or cosmetic fault classes ( $\alpha > .05$  for both). A distribution of the faults observed according to this classification scheme appears in Figure 17.

#### 4.3.3. Observable Fault Classification

Figure 18 displays the average percentage (with standard deviations) of faults from each class that were observable from the test data submitted, yet were not reported by the tester.<sup>5</sup> The two on-line techniques of functional and structural testing were not

---

<sup>5</sup> The standard deviations presented in the figure are high because of the several instances in which all observable faults were reported.



different in any of the faults classes ( $\alpha > .05$ ). Note that there was only one fault in the cosmetic class.

#### 4.3.4. Summary of Characterization of Faults Detected

Summarizing the major results of the comparison of classes of faults detected: 1) code reading and functional testing both detected more omission faults and initialization faults than did structural testing; 2) code reading detected more interface faults than did the other methods; 3) functional testing detected more control faults than did the other methods; 4) code reading detected more computation faults than did structural testing; and 5) the on-line techniques of functional and structural testing were not different in any classes of faults observable but not reported.

### 5. Conclusions

This study compares the strategies of code reading by stepwise abstraction, functional testing using equivalence class partitioning and boundary value analysis, and structural testing using 100% statement coverage. The study evaluates the techniques across three data sets in three different aspects of software testing: fault detection effectiveness, fault detection cost, and classes of faults detected. Each of the three testing techniques showed merit in this evaluation. The investigation is intended to compare the different testing strategies in representative testing situations, using programmers with a wide range of experience, different software types, and common software faults.

The major results of this study are 1) with the professional programmers, code reading detected more software faults and had a higher fault detection rate than did functional or structural testing, while functional testing detected more faults than did structural testing, but functional and structural testing were not different in fault detection rate; 2) in one UoM subject group, code reading and functional testing were not different in faults found, but were both superior to structural testing, while in the other UoM subject group there was no difference among the techniques; 3) with the UoM subjects, the three techniques were not different in fault detection rate; 4) number of faults observed, fault detection rate, and total effort in detection depended on the type of software tested; 5) code reading detected more interface faults than did the other methods;

6) functional testing detected more control faults than did the other methods; and 7) when asked to estimate the percentage of faults detected, code readers gave the most accurate estimates while functional testers gave the least accurate estimates.

The results suggest that code reading by stepwise abstraction (a nonexecution-based method) is at least as effective as on-line functional and structural testing in terms of number and cost of faults observed. They also suggest the inadequacy of using 100% statement coverage criteria for structural testing. Note that the professional programmers examined preferred the use of functional testing because they felt it was the most effective technique; their intuition, however, turned out to be incorrect.

In comparing the results to related studies, there are mixed conclusions. A prototype analysis done at the University of Maryland in the Fall of 1981 [Hwang 81] supported the belief that code reading by stepwise abstraction does as well as the computer-based methods, with each strategy having its own advantages. In the Myers experiment [Myers 78], the three techniques compared (functional testing, 3-person code reviews, control group) were equally effective. He also calculated that code reviews were less cost-effective than the computer-based testing approaches. The first observation is supported in one study phase here, but the other observation is not. A study conducted by Hetzel [Hetzel 76] compared functional testing, code reading, and "selective" testing (a composite of functional, structural, and reading techniques). He observed that functional and "selective" testing were equally effective, with code reading being inferior. As noted earlier, this is not supported by this analysis. The study described in this analysis examined the technique of code reading by stepwise abstraction, while both the Myers and Hetzel studies examined alternate approaches to off-line (nonexecution-based) review/reading.

A few remarks are appropriate about the comparison of the cost-effectiveness and phase-availability of these testing techniques. When examining the effort associated with a technique, both fault detection and fault isolation costs should be compared. The code readers have both detected and isolated a fault; they located it in the source code. Thus, the reading process condenses fault detection and isolation into one activity. Functional and structural testers have only detected a fault; they need to delve into the source code and expend additional effort in order to isolate the defect. Also, a

nonexecution-based reading process can be applied to any document produced during the development process (e.g., high-level design document, low-level design document, source code document). While functional and structural execution-based techniques may only be applied to documents that are executable (e.g., source code), which are usually available later in the development process.

Investigations related to this work include studies of fault classification [Weiss & Basili 85, Johnson, Draper & Soloway 83, Ostrand & Weyuker 83, Basili & Perricone 84] and Cleanroom software development [Selby, Basili & Baker 85]. In the Cleanroom software development approach, techniques such as code reading are used in the development of software completely off-line (i.e., without program execution). In the above study, systems developed using Cleanroom met system requirements more completely and had a higher percentage of successful operational test cases than did systems developed with a more traditional approach.

The empirical study presented is intended to advance the understanding of how various software testing strategies contribute to the software development process and to one another. The results given were calculated from a set of individuals applying the three techniques to unit-sized programs - the direct extrapolation of the findings to other testing environments is not implied. However, valuable insights into software testing have been gained. Further work applying these and other results to devise effective testing environments is underway.

## 6. Acknowledgement

The authors are grateful to the subjects from Computer Sciences Corporation, NASA Goddard, and the University of Maryland for their enthusiastic participation in the study.

## 7. Appendices

### 7.1. Appendix A. The Specifications for the Programs

#### *Program 1*

Given an input text of up to 80 characters consisting of words separated by blanks or new-line characters, the program formats it into a line-by-line form such that 1) each output line has a maximum of 30 characters, 2) a word in the input text is placed on a single output line, and 3) each output line is filled with as many words as possible.

The input text is a stream of characters, where the characters are categorized as either break or nonbreak characters. A break character is a blank, a new-line character (&), or an end-of-text character (/). New-line characters have no special significance; they are treated as blanks by the program. The characters & and / should not appear in the output.

A word is defined as a nonempty sequence of nonbreak characters. A break is a sequence of one or more break characters and is reduced to a single blank character or start of a new line in the output.

When the program is invoked, the user types the input line, followed by a / (end-of-text) and a carriage return. The program then echos the text input and formats it on the terminal.

If the input text contains a word that is too long to fit on a single output line, an error message is typed and the program terminates. If the end-of-text character is missing, an error message is issued and the program awaits the input of properly terminated line of text.

#### *Program 2*

Given ordered pairs (x,y) of either positive or negative integers as input, the program plots them on a grid with a horizontal x-axis and a vertical y-axis which are appropriately labeled. A plotted point on the grid should appear as an asterisk (\*).

The vertical and horizontal scaling is handled as follows. If the maximum absolute value of any y-value is less than or equal to twenty (20), the scale for vertical spacing will be one line per integral unit (e.g., the point (3,6) should be plotted on the sixth line; two lines above the point (3,4)). Note that the origin (point (0,0)) would correspond to an asterisk at the the intersection of the axes (the x-axis is referred to as the 0th line). If the maximum absolute value of any x-value is less than or equal to thirty (30), the scale for horizontal spacing will be one space per integral unit (e.g., the point (4,5) should be plotted four spaces to the right of the y-axis; two spaces to the right of (2,5)). However, if the maximum absolute value of any y-value is greater than twenty (20), the scale for vertical spacing will be one line per every (max abs of yval)/20 rounded-up. (e.g., If the maximum absolute value of any y-value to be plotted is 66, the vertical line spacing will be a line for every four (4) integral units. In such a data set, points with y-values greater than or equal to eight and less than twelve will show up as asterisks in the second line, points with y-values greater than or equal to twelve and less than six-

teen will show up as asterisks in the third line, etc. Continuing the example, the point (3,15) should be plotted on the third line; two lines above the point (3,5).) Horizontal scaling is handled analogously.

If two or more of the points to be plotted would show up as the same asterisk in the grid (like the points (9,13) and (9,15) in the above example), a number '2' (or whatever number is appropriate) should be printed instead of the asterisk. Points whose asterisks will lie on a axis or grid marker should show up in place of the marker.

### *Program 3*

A list is defined to be an ordered collection of integer elements which may have elements annexed and deleted at either end, but not in the middle. The operations that need to be available are ADDFIRST, ADDLAST, DELETFIRST, DELETELAST, FIRST, ISEMPY, LISTLENGTH, REVERSE, and NEWLIST. Each operation is described in detail below. The lists are to contain up to a maximum of five (5) elements. If an element is added to the front of a "full" list (one containing five elements already), the element at the back of the list is to be discarded. Elements to be added to the back of a full list are discarded. Requests to delete elements from empty lists result in an empty list, and requests for the first element of an empty list results in zero (0) being returned. The detailed operation descriptions are as below:

ADDFIRST(LIST L, INTEGER I)

Returns the list L with I as its first element followed by all the elements of L. If L is "full" to begin with, L's last element is lost.

ADDLAST(LIST L, INTEGER I)

Returns the list with all of the elements of L followed by I. If L is full to begin with, L is returned (i.e., I is ignored).

DELETFIRST(LIST L)

Returns the list containing all but the first element of L. If L is empty, then an empty list is returned.

DELETELAST(LIST L)

Returns the list containing all but the last element of L. If L is empty, then an empty list is returned.

FIRST(LIST L)

Returns the first element in L. If L is empty, then it returns zero (0).

ISEMPY(LIST L)

Returns one (1) if L is empty, zero (0) otherwise.

LISTLENGTH(LIST L)

Returns the number of elements in L. An empty list has zero (0) elements.

NEWLIST(LIST L)

Returns an empty list.

REVERSE(LIST L)

Returns a list containing the elements of L in reverse order.

#### Program 4

(Note that a 'file' is the same thing as an IBM 'dataset'.)

The program maintains a database of bibliographic references. It first reads a master file of current references, then reads a file of reference updates, merges the two, and produces an updated master file and a cross reference table of keywords.

The first input file, the master, contains records of 74 characters with the following format:

column	comment
--------	---------

---

1 - 3	each reference has a unique reference key
4 - 14	author of publication
15 - 72	title of publication
73 - 74	year issued

The key should be a three (3) character unique identifier consisting of letters between A-Z. The next input file, the update file, contains records of 75 characters in length. The only difference from a master file record is that an update record has either an 'A' (capital A meaning add) or a 'R' (capital R meaning replace) in column 75. Both the master and update files are expected to be already sorted alphabetically by reference key when read into the program. Update records with action replace are substituted for the matching key record in the master file. Records with action add are added to the master file at the appropriate location so that the file remains sorted on the key field. For example, a valid update record to be read would be (including a numbered line just for reference)

12345678901234567890123456789012345678901234567890123456789012345

BITbaker	an introduction to program testing	83A
----------	------------------------------------	-----

The program should produce two pieces of output. It should first print the sorted list of records in the updated master file in the same format as the original master file. It should then print a keyword cross reference list. All words greater than three characters in a publication's title are keywords. These keywords are listed alphabetically followed by the key fields from the applicable updated master file entries. For example, if the updated master file contained two records,

ABCKermitt	introduction to software testing	82
DDXJones	the realities of software management	81

then the keywords are introduction, testing, realities, software, and management. The cross reference list should look like

introduction
ABC
management
DDX
realities

DDX  
software  
ABC  
DDX  
testing  
ABC

Some possible error conditions that could arise and the subsequent actions include the following. The master and update files should be checked for sequence, and if a record out of sequence is found, a message similar to 'key ABC out of sequence' should appear and the record should be discarded. If an update record indicates replace and the matching key can not be found, a message similar to 'update key ABC not found' should appear and the update record should be ignored. If an update record indicates add and a matching key is found, something like 'key ABC already in file' should appear and the record should be ignored. (End of specification.)

## 7.2. Appendix B. The Source Code for the Programs

### *Program 1*

```
001: C NOTE THAT YOU DO NOT NEED TO VERIFY THE FUNCTION 'MATCH'.
002: C IT IS DESCRIBED THE FIRST TIME IT IS USED, AND ITS SOURCE CODE
003: C IS INCLUDED AT THE END FOR COMPLETENESS.
004: C
005: C NOTE THAT FORMAT STATEMENTS FOR WRITE STATEMENTS INCLUDE
      A LEADING
006: C AND REQUIRED ' ' FOR CARRIAGE CONTROL
007:
008: C VARIABLE USED IN FIRST, BUT NEEDS TO BE INITIALIZED
009:     INTEGER MOREIN
010:
011: C STORAGE USED BY GCHAR
012:     INTEGER BCOUNT
013:     CHARACTER*1 GBUFER(80)
014:     CHARACTER*80 GBUF
015: C GBUFER AND GBUFSTR ARE EQUIVALENCED
016:
017: C STORAGE USED BY PCHAR
018:     INTEGER I
019:     CHARACTER*1 OUTLIN(31)
020: C OUTLIN AND OUTLINST ARE EQUIVALENCED
021:
022:     CHARACTER*1 GCHAR
023:
024: C CONSTANT USED THROUGHOUT THE PROGRAM
025:     CHARACTER*1 EOTEXT, BLANK, LINEFD
026:     INTEGER MAXPOS
027:
028:     COMMON /ALL/ MOREIN, BCOUNT, I, MAXPOS, OUTLIN,
029:     X EOTEXT, BLANK, LINEFD, GBUFER, GBUF
030:
031:     DATA EOTEXT, BLANK, LINEFD, MAXPOS / '/', ' ', '&', 31 /
```

```

032:
033:
034:     CALL FIRST
035:     END
036:
037:
038:     SUBROUTINE FIRST
039:     INTEGER K, FILL, BUFPOS
040:     CHARACTER*1 CW
041:     CHARACTER*1 BUFFER(31)
042:
043:     INTEGER MOREIN, BCOUNT, I, MAXPOS
044:     CHARACTER*1 OUTLIN(31), GCHAR, EOTEXT, BLANK, LINEFD,
045: X      GBUFER(80)
046:     CHARACTER*80 GBUF
047:
048:     COMMON /ALL/ MOREIN, BCOUNT, I, MAXPOS, OUTLIN,
049: X      EOTEXT, BLANK, LINEFD, GBUFER, GBUF
050:
051:     BUFPOS = 0
052:     FILL = 0
053:     CW = ' '
054:
055:     MOREIN = 1
056:
057:     I = 1
058:     K = 1
059:     DOWHILE (K .LE. MAXPOS)
060:         OUTLIN(K) = ' '
061:         K = K + 1
062:     ENDDO
063:
064:     BCOUNT = 1
065:     K = 1
066:     DOWHILE (K .LE. 80)
067:         GBUFER(K) = 'Z'
068:         K = K + 1
069:     ENDDO
070:
071:     DOWHILE (MOREIN)
072:         CW = GCHAR()
073:         IF ((CW .EQ. BLANK) .OR. (CW .EQ. LINEFD) .OR.
074: X      (CW .EQ. EOTEXT)) THEN
075:             IF (CW .EQ. EOTEXT) THEN
076:                 MOREIN = 0
077:             ENDIF
078:             IF ((FILL+1+BUFPOS) .LE. MAXPOS) THEN
079:                 CALL PCHAR(BLANK)
080:                 FILL = FILL + 1
081:             ELSE
082:                 CALL PCHAR(LINEFD)
083:                 FILL = 0
084:             ENDIF
085:             K = 1
086:             DOWHILE (K .LE. BUFPOS)

```



```

087:          CALL PCHAR(BUFFER(K))
088:          K = K + 1
089:          ENDDO
090:          FILL = FILL + BUFPOS
091:          BUFPOS = 0
092:          ELSE
093:            IF (BUFPOS .EQ. MAXPOS) THEN
094:              WRITE(6,10)
095: 10          FORMAT(' ','***WORD TO LONG***')
096:              MOREIN = 1
097:            ELSE
098:              BUFPOS = BUFPOS + 1
099:              BUFFER(BUFPOS) = CW
100:            ENDIF
101:          ENDIF
102:        ENDDO
103:        CALL PCHAR(LINEFD)
104:        END
105:
106:
107:        CHARACTER*1 FUNCTION GCHAR()
108:        INTEGER MATCH
109:        CHARACTER*80 GBUFSTR
110:
111:        INTEGER MOREIN, BCOUNT, I, MAXPOS
112:        CHARACTER*1 OUTLIN(31), EOTEXT, BLANK, LINEFD,
113:  X      GBUFER(80)
114:        CHARACTER*80 GBUF
115:        COMMON /ALL/ MOREIN, BCOUNT, I, MAXPOS, OUTLIN,
116:  X      EOTEXT, BLANK, LINEFD, GBUFER, GBUF
117:
118:        EQUIVALENCE (GBUFSTR,GBUFER)
119:
120:        IF (GBUFER(1) .EQ. 'Z') THEN
121:          READ(5,20) GBUF
122: 20          FORMAT(A80)
123:  C
124:  C MATCH(CARRAY,C) RETURNS 1 IF CHARACTER C IS IN
        CHARACTER ARRAY
125:  C CARRAY, RETURNS 0 OTHERWISE. ARSIZE IS THE SIZE OF CARRAY.
126:  C
127:          IF (MATCH(GBUF,EOTEXT) .EQ. 0) THEN
128:            WRITE(6,30)
129: 30          FORMAT(' ','***NO END OF TEXT MARK***')
130:            GBUFER(2) = EOTEXT
131:          ELSE
132:  C          GBUFER(1) = GBUF
133:            GBUFSTR = GBUF
134:          ENDIF
135:        ENDIF
136:        GCHAR = GBUFER(BCOUNT)
137:        BCOUNT = BCOUNT + 1
138:        END
139:
140:

```

```

141: SUBROUTINE PCHAR (C)
142: CHARACTER*1 C
143: CHARACTER*31 SOUT, OUTLINST
144: INTEGER K
145:
146: INTEGER MOREIN, BCOUNT, I, MAXPOS
147: CHARACTER*1 OUTLIN(31), GCHAR, EOTEXT, BLANK, LINEFD,
148: X GBUFER(80)
149: CHARACTER*80 GBUF
150: COMMON /ALL/ MOREIN, BCOUNT, I, MAXPOS, OUTLIN,
151: X EOTEXT, BLANK, LINEFD, GBUFER, GBUF
152:
153: EQUIVALENCE (OUTLINST,OUTLIN)
154:
155: IF (C .EQ. LINEFD) THEN
156:     SOUT = OUTLINST
157:     WRITE(6,40) SOUT
158: 40    FORMAT(' ',A31)
159:     K = 1
160:     DOWHILE (K .LE. MAXPOS)
161:         OUTLIN(K) = ' '
162:         K = K + 1
163:     ENDDO
164:     I = 1
165: ELSE
166:     OUTLIN(I) = C
167:     I = I + 1
168: ENDIF
169: END

```

### Program 2

```

1: INT WIDTH = 30,
2:   HEIGHT = 20,
3:   GRIDWD = 61,
4:   LARGENUM = 100000000
5: STRING TICKS[61] =
6: ' |- |- |- |- |- |- |- |- |- |- |- |- |- |- |- '
7:
8:
9: PROC SORT (INT ARRAY KEYBUF, INT ARRAY FREEBUF, INT N)
10:
11: INT I, MAXP
12: INT ARRAY SRTKEYB(100), SRTFREEB(100)
13:
14: I := 0
15: WHILE I < N DO
16:     SRTKEYB(I) := KEYBUF(I)
17:     SRTFREEB(I) := FREEBUF(I)
18:     I := I + 1
19: END
20:
21: I := N
22: WHILE I > 0 DO
23:     MAXP := MAXELE(SRTKEYB,I)

```

```

24:   KEYBUF(N-I) := SRTKEYB(MAXP)
25:   FREEBUF(N-I) := SRTFREEB(MAXP)
26:   CALL REMOVE(SRTKEYB,MAXP,I)
27:   CALL REMOVE(SRTFREEB,MAXP,I)
28:   I := I - 1
29:   END
30:
31:
32:
33: INT FUNC MAXELE (INT ARRAY BUF, INT N)
34:
35:   INT I, MAXPTR, MAX
36:
37:   MAXPTR := -1
38:   MAX := -LARGENUM
39:   I := 0
40:   WHILE I < N DO
41:     IF BUF(I) > MAX
42:       THEN
43:         MAX := BUF(I)
44:         MAXPTR := I
45:       END
46:     I := I + 1
47:   END
48:   RETURN(MAXPTR)
49:
50:
51:
52: INT FUNC MINELE (INT ARRAY BUF, INT N)
53:
54:   INT I, MINPTR, MIN
55:
56:   MINPTR := -1
57:   MIN := LARGENUM
58:   I := 0
59:   WHILE I < N DO
60:     IF BUF(I) < MIN
61:       THEN
62:         MIN := BUF(I)
63:         MINPTR := I
64:       END
65:     I := I + 1
66:   END
67:   RETURN(MINPTR)
68:
69:
70:
71: PROC REMOVE (INT ARRAY BUF, INT PTR, INT N)
72:
73:   INT I
74:
75:   I := PTR
76:   WHILE I < N-1 DO
77:     BUF(I) := BUF(I+1)
78:     I := I + 1

```

```

79:     END
80:
81:
82:
83: INT FUNC ABS (INT VAL)
84:
85:     IF VAL < 0
86:     THEN
87:         RETURN(-VAL)
88:     ELSE
89:         RETURN(VAL)
90:     END
91:
92:
93:
94: INT FUNC SLASH (INT TOP, INT BOT)
95:
96:     INT RES
97:
98:     RES := TOP/BOT
99:     IF TOP <> RES*BOT .AND.
100:         (TOP > 0 .AND. BOT > 0 .OR. TOP < 0 .AND. BOT < 0)
101:     THEN RES := RES + 1
102:     END
103:     RETURN(RES)
104:
105: INT FUNC MOD (INT N, INT M)
106:
107:     INT VAL
108:
109:     VAL := N-N/M*M
110:     IF VAL < 0
111:     THEN
112:         VAL := VAL + M
113:     END
114:     RETURN (VAL)
115:
116:
117: PROC MAIN
118:
119:     CHAR ARRAY GRID(61)
120:     STRING STR(61)
121:     INT ARRAY XVAL(100), YVAL(100)
122:     INT I, J, NUMOBS, MAXY, MAXX, MINX, HORISP, VERTSP, VLINE
123:
124:     I := 0
125:     WHILE .NOT. EOI DO
126:         READ(XVAL(I),YVAL(I))
127:         I := I + 1
128:     END
129:     NUMOBS := I
130:
131:     CALL SORT(YVAL,XVAL,NUMOBS)
132:     MAXY := YVAL(0)
133:     VERTSP := SLASH(MAXY,HEIGHT)

```

```

134:
135:   MAXX := XVAL(MAXELE(XVAL,NUMOBS))
136:   MINX := XVAL(MINELE(XVAL,NUMOBS))
137:   IF ABS(MINX) > ABS(MAXX)
138:   THEN
139:     HORISP := SLASH(ABS(MINX),WIDTH)
140:   ELSE
141:     HORISP := SLASH(ABS(MAXX),WIDTH)
142:   END
143:
144:   STR := '                X AXIS'
145:   WRITE(STR,SKIP)
146:   I := 0
147:   VLINE := HEIGHT
148:   WHILE VLINE > 0 DO
149:
150:     J := 0
151:     IF MOD(VLINE,5) = 0
152:     THEN
153:       UNPACK(TICKS,GRID)
154:     ELSE
155:       WHILE J < GRIDWD DO
156:         GRID(J) := " "
157:         J := J + 1
158:       END
159:     END
160:
161:     VLINE := VLINE - 1
162:
163:     WHILE VLINE*VERTSP < YVAL(I) DO
164:       IF XVAL(I) >= 0
165:       THEN
166:         GRID(WIDTH + SLASH(XVAL(I),HORISP)) := "*"
167:       ELSE
168:         GRID(WIDTH - SLASH(-XVAL(I),HORISP)) := "*"
169:       END
170:       I := I + 1
171:     END
172:
173:     GRID(WIDTH) := "|"
174:     PACK(GRID,STR)
175:     WRITE(STR,SKIP)
176:   END
177:
178:   STR :=
179:   '-----'
180:   UNPACK(STR,GRID)
181:   WHILE 0 <= YVAL(I) .AND. I <= NUMOBS DO
182:     IF XVAL(I) >= 0
183:     THEN
184:       GRID(WIDTH + SLASH(XVAL(I),HORISP)) := "*"
185:     ELSE
186:       GRID(WIDTH - SLASH(-XVAL(I),HORISP)) := "*"
187:     END
188:     I := I + 1

```

```

188:      END
189:
190:      PACK(GRID,STR)
191:      WRITE(STR,SKIP)
192:      STR := '                Y AXIS'
193:      WRITE(STR,SKIP)
194:
195: START MAIN

```

### Program 3

```

001: C NOTE THAT YOU DO NOT NEED TO VERIFY THE FUNCTIONS
      DRIVER, GETARG,
002: C CHAREQ, CODE, AND PRINT. THEIR SOURCE CODE IS
      DESCRIBED AND
003: C INCLUDED AT THE END FOR COMPLETENESS.
004: C NOTE THAT FORMAT STATEMENTS FOR WRITE STATEMENTS
      INCLUDE A LEADING
005: C AND REQUIRED ' ' FOR CARRIAGE CONTROL
006: C
007:      INTEGER POOL(7), LSTEND
008:      INTEGER LISTSZ
009: C
010:      COMMON /ALL/ LISTSZ
011: C
012: C
013:      LISTSZ = 5
014:      CALL DRIVER (POOL, LSTEND)
015:      STOP
016:      END
017: C
018: C
019:      FUNCTION ADFRST (POOL, LSTEND, I)
020:      INTEGER ADFRST
021:      INTEGER POOL(7), LSTEND, I
022:      INTEGER LISTSZ
023:      COMMON /ALL/ LISTSZ
024: C
025:      INTEGER A
026: C
027:      IF (LSTEND .GT. LISTSZ) THEN
028:          LSTEND = LISTSZ - 1
029:      ENDIF
030:      LSTEND = LSTEND + 1
031:      A = LSTEND
032:      DOWHILE (A .GE. 1)
033:          POOL(A+1) = POOL(A)
034:          A = A - 1
035:      ENDDO
036: C
037:      POOL(1) = I
038:      ADFRST = LSTEND
039:      RETURN
040:      END
041: C

```

```

042: C
043: FUNCTION ADLAST (POOL, LSTEND, I)
044: INTEGER ADLAST
045: INTEGER POOL(7), LSTEND, I
046: INTEGER LISTSZ
047: COMMON /ALL/ LISTSZ
048: C
049: IF (LSTEND .LE. LISTSZ) THEN
050:     LSTEND = LSTEND + 1
051:     POOL(LSTEND) = I
052: ENDIF
053: ADLAST = LSTEND
054: RETURN
055: END
056: C
057: C
058: FUNCTION DELFST (POOL, LSTEND)
059: INTEGER DELFST
060: INTEGER POOL(7), LSTEND
061: INTEGER LISTSZ
062: COMMON /ALL/ LISTSZ
063: C
064: INTEGER A
065: IF (LSTEND .GT. 1) THEN
066:     A = 1
067:     LSTEND = LSTEND - 1
068:     DOWHILE (A .LE. LSTEND)
069:         POOL(A) = POOL(A+1)
070:         A = A + 1
071:     ENDDO
072: ENDIF
073: DELFST = LSTEND
074: RETURN
075: END
076: C
077: C
078: FUNCTION DELLST (LSTEND)
079: INTEGER DELLST
080: INTEGER LSTEND
081: C
082: IF (LSTEND .GE. 1) THEN
083:     LSTEND = LSTEND - 1
084: ENDIF
085: DELLST = LSTEND
086: RETURN
087: END
088: C
089: C
090: FUNCTION FIRST (POOL, LSTEND)
091: INTEGER FIRST
092: INTEGER POOL(7), LSTEND
093: C
094: IF (LSTEND .LE. 1) THEN
095:     FIRST = 0
096: ELSE

```

```

097:      FIRST = POOL(1)
098:      ENDIF
099:      RETURN
100:      END
101: C
102: C
103:      FUNCTION EMPTY (LSTEND)
104:      INTEGER EMPTY
105:      INTEGER LSTEND
106: C
107:      IF (LSTEND .LE. 1) THEN
108:          EMPTY = 1
109:      ELSE
110:          EMPTY = 0
111:      ENDIF
112:      RETURN
113:      END
114: C
115: C
116:      FUNCTION LSTLEN (LSTEND)
117:      INTEGER LSTLEN
118:      INTEGER LSTEND
119: C
120:      LSTLEN = LSTEND - 1
121:      RETURN
122:      END
123: C
124: C
125:      FUNCTION NEWLST (LSTEND)
126:      INTEGER NEWLST
127:      INTEGER LSTEND
128: C
129:      NEWLST = 0
130:      RETURN
131:      END
132: C
133: C
134:      SUBROUTINE REVERS (POOL, LSTEND)
135:      INTEGER POOL(7), LSTEND
136: C
137:      INTEGER I, N
138: C
139:      N = LSTEND
140:      I = 1
141:      DOWHILE (I .LE. N)
142:          POOL(I) = POOL(N)
143:          N = N - 1
144:          I = I + 1
145:      ENDDO
146:      RETURN
147:      END

```

*Program 4*

001: C NOTE THAT YOU DO NOT NEED TO VERIFY THE ROUTINES



```

DRIVER, STREQ, WORDEQ,
002: C  NXTSTR, ARRCPY, CHARPT, BEFORE, CHAREQ, AND WRDBEF.
      THEIR SOURCE
003: C  CODE IS DESCRIBED AND INCLUDED AT THE END FOR
      COMPLETENESS.
004: C  NOTE THAT FORMAT STATEMENTS FOR WRITE STATEMENTS
      INCLUDE A LEADING
005: C  AND REQUIRED ' ' FOR CARRIAGE CONTROL
006: C  THE SFORT LANGUAGE CONSTRUCT '.IF (EXPRESSION)' BEGINS
      A BLOCKED
007: C  IF-THEN[-ELSE] STATEMENT, AND IT IS EQUIVALENT TO
      THE F77
008: C  'IF (EXPRESSION) THEN'.
009: C
010:    CALL DRIVER
011:    STOP
012:    END
013: C
014: C
015:    SUBROUTINE MAINSB
016: C
017:    LOGICAL*1 U$KEY(3),U$AUTH(11),U$TITL(58),U$YEAR(2),U$ACTN(1)
018:    LOGICAL*1 M$KEY(3),M$AUTH(11),M$TITL(58),M$YEAR(2)
019:    LOGICAL*1 ZZZ(3), LASTUK(3), LASTMK(3)
020:    LOGICAL*1 STREQ, CHAREQ, BEFORE, CHARPT
021:    INTEGER I
022: C
023:    LOGICAL*1 WORD(500,12), REFKEY(1000,3)
024:    INTEGER NUMWDS, NUMREF, PTR(500), NEXT(1000)
025:    COMMON /WORDS/ WORD, REFKEY, NUMWDS, NUMREF, PTR, NEXT
026: C
027:    WRITE(6,290)
028: 290  FORMAT(' ', '  UPDATED LIST OF MASTER ENTRIES')
029:    DO 300 I = 1, 3
030:      LASTMK(I) = CHARPT(' ')
031:      LASTUK(I) = CHARPT(' ')
032:      ZZZ(I) = CHARPT('Z')
033: 300  CONTINUE
034: C
035:    NUMWDS = 0
036:    NUMREF = 0
037:    CALL GETNM(M$KEY,M$AUTH,M$TITL,M$YEAR,LASTMK)
038:    CALL GETNUP(U$KEY,U$AUTH,U$TITL,U$YEAR,U$ACTN,LASTUK)
039: C
040:    DOWHILE ((.NOT.(STREQ(M$KEY,ZZZ,3))) .OR.
041:  X      (.NOT.(STREQ(U$KEY,ZZZ,3))) )
042:      .IF (STREQ(U$KEY,M$KEY,3))
043:      .IF (.NOT.(CHAREQ(U$ACTN(1),'R'))))
044:        WRITE(6,100) U$KEY
045: 100  FORMAT(' ', 'KEY ',3A1,' IS ALREADY IN FILE')
046:      ENDIF
047:      CALL OUTPUT(U$KEY,U$AUTH,U$TITL,U$YEAR)
048:      CALL DICTUP(U$KEY,U$TITL,58)
049:      CALL GETNM(M$KEY,M$AUTH,M$TITL,M$YEAR,LASTMK)
050:      CALL GETNUP(U$KEY,U$AUTH,U$TITL,U$YEAR,U$ACTN,LASTUK)

```

```

051:      ENDIF
052: C
053:      .IF (BEFORE(M$KEY,3,U$KEY,3))
054:          CALL OUTPUT(M$KEY,M$AUTH,M$TITL,M$YEAR)
055:          CALL DICTUP(M$KEY,M$TITL,58)
056:          CALL GETNM(M$KEY,M$AUTH,M$TITL,M$YEAR,LASTMK)
057:      ENDIF
058: C
059:      .IF (BEFORE(U$KEY,3,M$KEY,3))
060:          .IF (CHAREQ(U$ACTN(1),R'))
061:              WRITE(6,110) U$KEY
062: 110      FORMAT(' ',UPDATE KEY ',3A1,' NOT FOUND')
063:          ENDIF
064:          CALL OUTPUT(U$KEY,U$AUTH,U$TITL,U$YEAR)
065:          CALL DICTUP(U$KEY,U$TITL,58)
066:          CALL GETNUP(U$KEY,U$AUTH,U$TITL,U$YEAR,U$ACTN,LASTUK)
067:      ENDIF
068: ENDDO
069: C
070: CALL SRTWDS
071: CALL PRTWDS
072: RETURN
073: END
074: C
075: C
076: SUBROUTINE GETNM(KEY,AUTH,TITL,YEAR,LASTMK)
077: LOGICAL*1 KEY(3),AUTH(11),TITL(58),YEAR(2),LASTMK(3)
078: C
079: LOGICAL*1 SEQ, INLINE(80)
080: LOGICAL*1 BEFORE, CHARPT, CHAREQ
081: LOGICAL*1 GO$M, GO$U
082: COMMON /DRIV/ GO$M, GO$U
083: C
084: SEQ = 1
085: DOWHILE (SEQ)
086:     .IF (GO$M)
087: C
088: C READ FROM THE MASTER FILE
089: C
090:     READ(10,200,END=299) INLINE
091:     ELSE
092: C
093: C SEE REMARK ABOUT THE CHARACTER '%' LATER IN THE ROUTINE.
094: C
095:     INLINE(1) = CHARPT('%')
096:     ENDIF
097: 200     FORMAT(80A1)
098:     DO 210 I = 1, 3
099:         KEY(I) = INLINE(I)
100: 210     CONTINUE
101:     DO 220 I = 1, 11
102:         AUTH(I) = INLINE(3+I)
103: 220     CONTINUE
104:     DO 230 I = 1, 58
105:         TITL(I) = INLINE(14+I)

```

```

106: 230    CONTINUE
107:        DO 240 I = 1, 2
108:        YEAR(I) = INLINE(72+I)
109: 240    CONTINUE
110: C
111: C A METHOD OF SPECIFYING END-OF-FILE IN A FILE IS TO PUT
    THE CHARACTER '%'
112: C AS THE FIRST CHARACTER ON A LINE. THE DRIVER USES THIS
    FOR MULTIPLE
113: C SETS OF INPUT CASES.
114: C
115:        .IF ((.NOT.(CHAREQ(KEY(1),'%')))) AND.
116: X      (BEFORE(KEY,3,LASTMK,3)) )
117:        WRITE(6,250) KEY
118: 250    FORMAT(' ',KEY ',3A1,' OUT OF SEQUENCE')
119:        ELSE
120:        CALL ARRCPY(KEY,LASTMK,3)
121:        SEQ = 0
122:        ENDIF
123:        .IF (CHAREQ(KEY(1),'%'))
124:        SEQ = 0
125:        DO 270 I = 1, 3
126:        KEY(I) = CHARPT('Z')
127: 270    CONTINUE
128:        ENDIF
129:        ENDDO
130:        RETURN
131: 299    CONTINUE
132:        GO$M = 0
133:        DO 260 I = 1, 3
134:        KEY(I) = CHARPT('Z')
135: 260    CONTINUE
136:        RETURN
137:        END
138: C
139: C
140: SUBROUTINE GETNUP(KEY,AUTH,TITL,YEAR,ACTN,LASTUK)
141: LOGICAL*1 KEY(3),AUTH(11),TITL(58),YEAR(2),ACTN(1),LASTUK(3)
142: C
143: LOGICAL*1 SEQ, INLINE(80)
144: LOGICAL*1 BEFORE, CHARPT, CHAREQ
145: LOGICAL*1 GO$M, GO$U
146: COMMON /DRIV/ GO$M, GO$U
147: C
148: SEQ = 1
149: DOWHILE (SEQ)
150: .IF (GO$U)
151: C
152: C READ FROM THE UPDATES FILE
153: C
154: READ(11,200,END=299) INLINE
155: ELSE
156: C
157: C SEE REMARK ABOUT THE CHARACTER '%' LATER IN THE ROUTINE.
158: C

```

```

159:      INLINE(1) = CHARPT('%')
160:      ENDIF
161: 200      FORMAT(80A1)
162:      DO 210 I = 1, 3
163:          KEY(I) = INLINE(I)
164: 210      CONTINUE
165:      DO 220 I = 1, 11
166:          AUTH(I) = INLINE(3+I)
167: 220      CONTINUE
168:      DO 230 I = 1, 58
169:          TITL(I) = INLINE(14+I)
170: 230      CONTINUE
171:      DO 240 I = 1, 2
172:          YEAR(I) = INLINE(72+I)
173: 240      CONTINUE
174:          ACTN(1) = INLINE(75)
175: C
176: C A METHOD OF SPECIFYING END-OF-FILE IN A FILE IS TO PUT
      THE CHARACTER '%'
177: C AS THE FIRST CHARACTER ON A LINE. THE DRIVER USES THIS
      FOR MULTIPLE
178: C SETS OF INPUT CASES.
179: C
180:      .IF ((.NOT.(CHAREQ(KEY(1),'%')) AND.
181:  X      (BEFORE(KEY,3,LASTUK,3)) )
182:          WRITE(6,250) KEY
183: 250      FORMAT(' ',KEY ',3A1,' OUT OF SEQUENCE')
184:      ELSE
185:          CALL ARRCPY(KEY,LASTUK,3)
186:          SEQ = 0
187:      ENDIF
188:      .IF (CHAREQ(KEY(1),'%'))
189:          SEQ = 0
190:          DO 270 I = 1, 3
191:              KEY(I) = CHARPT('Z')
192: 270      CONTINUE
193:      ENDIF
194:      ENDDO
195:      RETURN
196: 299 CONTINUE
197:      GO$U = 0
198:      DO 280 I = 1, 3
199:          KEY(I) = CHARPT('Z')
200: 280 CONTINUE
201:      RETURN
202:      END
203: C
204: C
205:      SUBROUTINE OUTPUT(KEY,AUTH,TITL,YEAR)
206:      LOGICAL*1 KEY(3), AUTH(11), TITL(58), YEAR(2)
207: C
208:      WRITE(6,200) KEY, AUTH, TITL, YEAR
209: 200      FORMAT(' ',3A1,11A1,58A1,2A1)
210:      RETURN
211:      END

```

```

212: C
213: C
214:   SUBROUTINE PRTWDS
215: C
216:   LOGICAL*1 WORD(500,12), REFKEY(1000,3)
217:   INTEGER NUMWDS, NUMREF, PTR(500), NEXT(1000)
218:   COMMON /WORDS/ WORD, REFKEY, NUMWDS, NUMREF, PTR, NEXT
219: C
220: C THE ABOVE GROUP OF DATA STRUCTURES SIMULATES A LINKED
    LIST.
221: C WORD(I,J) IS A KEYWORD - J RANGING FROM 1 TO 12
222: C REFKEY(PTR(I),K),K=1,3 IS THE FIRST 3 LETTER KEY THAT HAS
    AS A
223: C   KEYWORD WORD(I,J),J=1,12
224: C REFKEY(NEXT(PTR(I)),K),K=1,3 IS THE SECOND 3 LETTER KEY
    THAT HAS
225: C   AS A KEYWORD WORD(I,J),J=1,12
226: C REFKEY(NEXT(NEXT(PTR(I))),K),K=1,3 IS THE THIRD ... ETC.
227: C NEXT(J) IS EQUAL TO -1 WHEN THERE ARE NO MORE 3 LETTER
    KEYS FOR
228: C   THE PARTICULAR KEYWORD
229: C
230:   INTEGER I, J
231:   LOGICAL*1 FLAG
232: C
233:   WRITE(6,200)
234: 200   FORMAT(' ', ' KEYWORD REFERENCE LIST')
235:   DO 210 I = 1, NUMWDS
236:     FLAG = 1
237:     WRITE(6,220) (WORD(I,J),J=1,12)
238: 220   FORMAT(' ', '12A1)
239:     LAST = PTR(I)
240:     DOWHILE (FLAG)
241:       WRITE(6,230) (REFKEY(LAST,J),J=1,3)
242: 230   FORMAT(' ', '3A1)
243:       LAST = NEXT(LAST)
244:       .IF (LAST .EQ. -1)
245:         FLAG = 0
246:       ENDIF
247:     ENDDO
248: 210   CONTINUE
249:   RETURN
250:   END
251: C
252: C
253:   SUBROUTINE DICTUP(KEY,STR,STRLEN)
254:   LOGICAL*1 KEY(3), STR(120)
255:   INTEGER STRLEN
256: C
257:   LOGICAL*1 WDLEFT, FLAG, OKLEN, NEXTWD(120), WORDEQ
258:   INTEGER LPTR, NXTSTR, LEN, LAB, I, K
259: C
260:   LOGICAL*1 WORD(500,12), REFKEY(1000,3)
261:   INTEGER NUMWDS, NUMREF, PTR(500), NEXT(1000)
262:   COMMON /WORDS/ WORD, REFKEY, NUMWDS, NUMREF, PTR, NEXT

```

```

263: C
264: C THE ABOVE GROUP OF DATA STRUCTURES SIMULATES A
    LINKED LIST.
265: C WORD(I,J) IS A KEYWORD - J RANGING FROM 1 TO 12
266: C REFKEY(PTR(I),K),K=1,3 IS THE FIRST 3 LETTER KEY THAT HAS
    AS A
267: C KEYWORD WORD(I,J),J=1,12
268: C REFKEY(NEXT(PTR(I)),K),K=1,3 IS THE SECOND 3 LETTER KEY
    THAT HAS
269: C AS A KEYWORD WORD(I,J),J=1,12
270: C REFKEY(NEXT(NEXT(PTR(I))),K),K=1,3 IS THE THIRD ... ETC.
271: C NEXT(J) IS EQUAL TO -1 WHEN THERE ARE NO MORE 3 LETTER
    KEYS FOR
272: C THE PARTICULAR KEYWORD
273: C
274:   WDLEFT = 1
275:   LPTR = 1
276: C
277:   DOWHILE (WDLEFT)
278:     FLAG = 1
279:     OKLEN = 1
280:     LEN = NXTSTR(STR,STRLEN,LPTR,NEXTWD,120)
281:     .IF (LEN .EQ. 0)
282:       WDLEFT = 0
283:     ENDIF
284: C
285:     .IF (LEN .LE. 2)
286:       OKLEN = 0
287:     ENDIF
288: C
289:     .IF (OKLEN)
290:       I = 1
291:       DOWHILE ((I .LE. NUMWDS).AND. FLAG )
292:         .IF (WORDEQ(NEXTWD,I))
293:           LAB = I
294:           FLAG = 0
295:         ENDIF
296:         I = I + 1
297:       ENDDO
298:       .IF (FLAG)
299:         NUMWDS = NUMWDS + 1
300:         NUMREF = NUMREF + 1
301:         DO 300 K = 1, 12
302:           WORD(NUMWDS,K) = NEXTWD(K)
303: 300       CONTINUE
304:         PTR(NUMWDS) = NUMREF
305:         DO 310 K = 1, 3
306:           REFKEY(NUMREF,K) = KEY(K)
307: 310       CONTINUE
308:         NEXT(NUMREF) = -1
309:       ELSE
310:         NUMREF = NUMREF + 1
311:         DO 320 K = 1, 3
312:           REFKEY(NUMREF,K) = KEY(K)
313: 320       CONTINUE

```

```

314:         NEXT(NUMREF) = PTR(LAB)
315:         PTR(LAB) = NUMREF
316:     ENDIF
317: ENDIF
318: ENDDO
319: C
320: RETURN
321: END
322: C
323: C
324: SUBROUTINE SRTWDS
325: C
326: LOGICAL*1 WORD(500,12), REFKEY(1000,3)
327: INTEGER NUMWDS, NUMREF, PTR(500), NEXT(1000)
328: COMMON /WORDS/ WORD, REFKEY, NUMWDS, NUMREF, PTR, NEXT
329: C
330: C THE ABOVE GROUP OF DATA STRUCTURES SIMULATES A
    LINKED LIST.
331: C WORD(I,J) IS A KEYWORD - J RANGING FROM 1 TO 12
332: C REFKEY(PTR(I),K),K=1,3 IS THE FIRST 3 LETTER KEY THAT HAS
    AS A
333: C KEYWORD WORD(I,J),J=1,12
334: C REFKEY(NEXT(PTR(I)),K),K=1,3 IS THE SECOND 3 LETTER KEY
    THAT HAS
335: C AS A KEYWORD WORD(I,J),J=1,12
336: C REFKEY(NEXT(NEXT(PTR(I))),K),K=1,3 IS THE THIRD ... ETC.
337: C NEXT(J) IS EQUAL TO -1 WHEN THERE ARE NO MORE 3 LETTER
    KEYS FOR
338: C THE PARTICULAR KEYWORD
339: C
340: INTEGER I, J, K, LAB, LOWERB, UPPERB
341: LOGICAL*1 WRDBEF, NEXTWD(12)
342: C
343: UPPERB = NUMWDS - 1
344: DO 400 I = 1, UPPERB
345:     LOWERB = I + 1
346:     DO 410 J = LOWERB, NUMWDS
347:         .IF (WRDBEF(J,I))
348:             DO 300 K = 1, 12
349:                 NEXTWD(K) = WORD(I,K)
350: 300 CONTINUE
351:         LAB = PTR(I)
352:         DO 310 K = 1, 12
353:             WORD(I,K) = WORD(J,K)
354: 310 CONTINUE
355:         PTR(I) = PTR(J)
356:         DO 320 K = 1, 12
357:             WORD(J,K) = NEXTWD(K)
358: 320 CONTINUE
359:         PTR(J) = LAB
360:     ENDIF
361: 410 CONTINUE
362: 400 CONTINUE
363: C
364: RETURN

```

365: END



## 8. References

[Basili & Turner 76]

V. R. Basili and A. J. Turner, *SIMPL-T: A Structured Programming Language*, Paladin House Publishers, Geneva, IL, 1976.

[Basili & Perricone 84]

V. R. Basili and B. T. Perricone, Software Errors and Complexity: An Empirical Investigation, *Communications of the ACM* 27, 1, pp. 42-52, Jan. 1984.

[Basili & Selby 84]

V. R. Basili and R. W. Selby, Jr., Data Collection and Analysis in Software Research and Management, *Proceedings of the American Statistical Association and Biometric Society Joint Statistical Meetings*, Philadelphia, PA, August 13-16, 1984.

[Basili & Weiss 84]

V. R. Basili and D. M. Weiss, A Methodology for Collecting Valid Software Engineering Data\*, *Trans. Software Engr.* SE-10, 6, pp. 728-738, Nov. 1984.

[Box, Hunter, & Hunter 78]

G. E. P. Box, W. G. Hunter, and J. S. Hunter, *Statistics for Experimenters*, John Wiley & Sons, New York, 1978.

[Callilau & Rubin 79]

R. Callilau and F. Rubin, ACM Forum: On a Controlled Experiment in Program Testing, *Communications of the ACM* 22, pp. 687-8, Dec. 1979.

[Church 84]

V. Church, Benchmark Statistics for the VAX 11/780 and the IBM 4341, Computer Sciences Corporation, Silver Spring, MD, Internal Memo, 1984.

[Cochran & Cox 50]

W. G. Cochran and G. M. Cox, *Experimental Designs*, John Wiley & Sons, New York, 1950.

[Fagan 76]

M. E. Fagan, Design and Code Inspections to Reduce Errors in Program Development, *IBM Sys. J.* 15, 3, pp. 182-211, 1976.

[Foster 80]

K. A. Foster, Error Sensitive Test Cases, *IEEE Trans. Software Engr.* SE-6, 3, pp. 258-264, 1980.

[Gloss-Soler 79]

S. A. Gloss-Soler, The DACS Glossary: A Bibliography of Software Engineering Terms, Data & Analysis Center for Software, Griffiss Air Force Base, NY 13441, Rep. GLOS-1, Oct. 1979.

[Goodenough & Gerhart 75]

J. B. Goodenough and S. L. Gerhart, Toward a Theory of Test Data Selection, *IEEE Trans. Software Engr.*, pp. 156-173, June 1975.

[Hetzel 76]

W. C. Hetzel, An Experimental Analysis of Program Verification Methods, Ph.D. Thesis, Univ. of North Carolina, Chapel Hill, 1976.

[Howden 78]

W. E. Howden, Algebraic Program Testing, *Acta Informatica* 10, 1978.

[Howden 80]

W. E. Howden, Functional Program Testing, *IEEE Trans. Software Engr.* SE-6, pp. 162-169, Mar. 1980.

[Howden 81]

W. E. Howden, A Survey of Dynamic Analysis Methods, pp. 209-231 in *Tutorial: Software Testing & Validation Techniques, 2nd Ed.*, ed. E. Miller and W. E. Howden, 1981.

[Hwang 81]

S-S. V. Hwang, An Empirical Study in Functional Testing, Structural Testing, and Code Reading/Inspection\*, Dept. Com. Sci., Univ. of Maryland, College Park, Scholarly Paper 362, Dec. 1981.

[IEEE 83]

IEEE, IEEE Standard Glossary of Software Engineering Terminology, Rep. IEEE-STD-729-1983, IEEE, 342 E. 47th St, New York, 1983.

[Jensen & Wirth 74]

K. Jensen and N. Wirth, *PASCAL User Manual and Report, 2nd Ed.*, Springer-Verlag, New York, 1974.

[Johnson, Draper & Soloway 83]

W. L. Johnson, S. Draper, and E. Soloway, An Effective Bug Classification Scheme Must Take the Programmer into Account, *Proc. Workshop High-Level Debugging*, Palo Alto, CA, 1983.

[Linger, Mills & Witt 79]

R. C. Linger, H. D. Mills, and B. I. Witt, *Structured Programming: Theory and Practice*, Addison-Wesley, Reading, MA, 1979.

[McMullin & Gannon 80]

P. R. McMullin and J. D. Gannon, Evaluating a Data Abstraction Testing System Based on Formal Specifications, Dept. Com. Sci., Univ. of Maryland, College Park, Tech. Rep. TR-993, Dec. 1980.

[Mills 72]

H. D. Mills, Mathematical Foundations for Structural Programming, IBM Report FSL 72-6021, 1972.

[Mills 75]

H. D. Mills, How to Write Correct Programs and Know It, *Int. Conf. on Reliable Software*, Los Angeles, pp. 363-370, 1975.

[Myers 78]

G. J. Myers, A Controlled Experiment in Program Testing and Code Walkthroughs/Inspections, *Communications of the ACM*, pp. 760-768, Sept. 1978.

[Myers 79]

G. J. Myers, *The Art of Software Testing*, John Wiley & Sons, New York, 1979.

[Naur 69]

P. Naur, Programming by Action Clusters, *BIT* 9, 3, pp. 250-258, 1969.

[Ostrand & Weyuker 83]

T. J. Ostrand and E. J. Weyuker, Collecting and Categorizing Software Error Data in an Industrial Environment, Dept. Com. Sci., Courant Inst. Math. Sci., New York Univ., NY, Tech. Rep. 47, August 1982 (Revised May 1983).

[Selby 83]

R. W. Selby, Jr., An Empirical Study Comparing Software Testing Techniques, *Sixth Minnowbrook Workshop on Software Performance Evaluation*, Blue Mountain Lake, NY, July 19-22, 1983.

[Selby 84]

R. W. Selby, Jr., Evaluating Software Testing Strategies, *Proc. of the Ninth Annual Software Engineering Workshop*, NASA/GSFC, Greenbelt, MD, Nov. 1984.

[Selby, Basill & Baker 85]

R. W. Selby, Jr., V. R. Basill, and F. T. Baker, CLEANROOM Software Development: An Empirical Evaluation, Dept. Com. Sci., Univ. Maryland, College Park, Tech. Rep. TR-1415, February 1985. (submitted to the *IEEE Trans. Software Engr.*)

[Stuckl 77]

L. G. Stuckl, New Directions in Automated Tools for Improving Software Quality, in *Current Trends in Programming Methodology*, ed. R. T. Yeh, Prentice Hall, Englewood Cliffs, NJ, 1977.

[Valdes & Goel 83]

P. M. Valdes and A. L. Goel, An Error-Specific Approach to Testing, *Proc. Eight Ann. Software Engr. Workshop*, NASA/GSFC, Greenbelt, MD, Nov. 1983.

[Weiss & Basill 85]

D. M. Weiss and V. R. Basill, Evaluating Software Development by Analysis of Changes: Some Data from the Software Engineering Laboratory, *IEEE Trans. Software Engr.* SE-11, 2, pp. 157-168, February 1985.

Figure 1. Capabilities of the testing methods.			
	code reading	functional testing	structural testing
view program specification	X	X	X
view source code	X		X
execute program		X	X

Figure 2. Structure of goals/subgoals/questions for testing experiment.

#### I. Fault detection effectiveness

A. For programmers doing unit testing, which of the testing techniques (code reading, functional testing, or structural testing) detects the most faults in programs?

1. Which of the techniques detects the greatest percentage of faults in the programs (the programs each contain a different number of faults)?
2. Which of the techniques exposes the greatest number (or percentage) of program faults (faults that are observable but not necessarily reported)?

B. Is the number of faults observed dependent on software type?

C. Is the number of faults observed dependent on the expertise level of the person testing?

#### II. Fault detection cost

A. For programmers doing unit testing, which of the testing techniques (code reading, functional testing, or structural testing) detects the faults at the highest rate (#faults/effort)?

B. Is the fault detection rate dependent on software type?

C. Is the fault detection rate dependent on the expertise level of the person testing?

### III. Classes of faults observed

A. For programmers doing unit testing, do the methods tend to capture different classes of faults?

B. What classes of faults are observable but go unreported?

Figure 3. Expertise levels of subjects.

Level of Expertise	Phase			total
	1 (Unlv. Md)	2 (Unlv. Md)	3 (NASA/CSC)	
Advanced	0	0	8	8
Intermediate	9	4	11	24
Junlor	20	9	13	42
total	29	13	32	74

Figure 4. The programs tested.

program	source lines	executable statments	cyclomatic complexity	#routines	#faults
$P_1$ - text formatter	169	55	18	3	9
$P_2$ - mathematical plotting	145	95	32	8	6
$P_3$ - numeric data abstraction	147	48	18	9	7
$P_4$ - database maintainer	365	144	57	7	12

Figure 5. Programs tested in each phase of the analysis.

Program	Phase		
	1 (Unlv. Md)	2 (Unlv. Md)	3 (NASA/CSC)
$P_1$ - text formatter	X	X	X
$P_2$ - mathematical plotting	X	X	
$P_3$ - numeric data abstraction	X		X
$P_4$ - database maintainer		X	X

Figure 6. Distribution of faults in the programs.

	Omission	Commission	Total
Initialization	0	2	2
Computation	4	4	8
Control	2	5	7
Interface	2	11	13
Data	2	1	3
Cosmetic	0	1	1
Total	10	24	34

Figure 7. Fault classification and manifestation.

Fault	Program	Omission/ Commission	Class	Description
a	P1	omission	control	a blank is printed before the first word on the first line unless the first word is 30 characters long; in the latter case, a blank line is printed before the first word
b	P1	commission	initialization	the character & (not \$) is the new-line character
c	P1	commission	initialization	the line size is 31 characters (not 30); this fault causes the references to the number 30 in the other faults to be actually the number 31
d	P1	commission	interface	since the program pads an empty input buffer with the character "z," it ignores a valid input line that has a "z" as a first character
e	P1	omission	control	successive break characters are not condensed in the output

f	P1	commisslon	cosmetic	spelling mistake in the error message "*** word to long ***"
g	P1	commisslon	computation	after detecting a word in the input longer than 30 characters, the message "*** word to long ***" is printed once for every character over 30, and the processing of the text does not terminate
h	P1	omisslon	interface	after detecting a word in the input longer than 30 characters, the program prints whatever is residing in its output buffer
i	P1	commisslon	control	after detecting an input line without an end-of-text character, the program erroneously increments its buffer pointer and replaces the first character of the next input line with a "z"
j	P3	commisslon	interface	routine FIRST returns zero (0) when the list has one element
k	P3	commlsslon	interface	routine ISEMPY returns true (1) when the list has one element
l	P3	commisslon	interface	routine DELETEFIRST can not delete the first list element when the list has only one element
m	P3	commisslon	interface	routine LISTLENGTH returns one less than than the actual length of the list
n	P3	commisslon	interface	routine ADDFIRST can add more than the specified five elements to the list
o	P3	commisslon	interface	routine ADDLAST can add more than the specified five elements to the list
p	P3	omisslon	computation	routine REVERSE does not reverse the list properly when the list has more than one element
q	P4	commlsslon	computation	words greater than or equal to three characters (not strictly greater than) are treated as cross reference keywords
r	P4	commisslon	interface	since the program uses the key "ZZZ" as an end-of-input sentinel, it does not process a valid record with key "ZZZ" and ignores any following records
s	P4	commlsslon	control	update action add with the error condition "key already in the master file" replaces the existing record; the update record is not ignored



t	P4	commission	control	update action replace with the error condition "key not found in the master file" adds the record; the update record is not ignored
u	P4	omission	data	the number of references and number of words in the dictionary are not checked for overflow
v	P4	omission	computation	two or more update transactions for the same master record give incorrect results
w	P4	commission	interface	keywords longer than 12 characters are truncated and not distinguished
x	P4	commission	control	an update record with column 80 neither an add action "A" nor replace action "R" acts like an add transaction
y	P4	commission	interface	keyword indices appear in reverse alphabetical order
z	P4	omission	interface	no check is made for unique keys in the master file
A	P4	commission	interface	punctuation is made a part of the keyword
B	P4	omission	data	words appearing twice in a title get two cross reference entries
C	P2	commission	computation	the x and y axes are mislabeled
D	P2	omission	computation	points with negative y-values are not processed and do not appear on the graph
E	P2	commission	control	the origin (0,0) appears on the graph regardless of whether it is an input point
F	P2	commission	data	no points can appear on the vertical axis
G	P2	commission	computation	the vertical and horizontal scaling for the pixels are calculated incorrectly, causing some points not to appear in the proper pixel
H	P2	omission	computation	when more than one point would appear in a given pixel, only an asterisk (*) appears, not an appropriate integer

Figure 8. Fractional Factorial Design.				
		Code Reading	Functional Testing	Structural Testing
		$P_1 P_3 P_4$	$P_1 P_3 P_4$	$P_1 P_3 P_4$
Advanced Subjects	$S_1$	—X	—X—	X—
	$S_2$	—X—	X—	—X
	.		...	
	$S_8$	X—	—X	—X—
Inter- medate Subjects	$S_9$	—X—	X—	—X
	$S_{10}$	—X	—X—	X—
	.		...	
	$S_{19}$	X—	—X	—X—
Junior Subjects	$S_{20}$	—X—	X—	—X
	$S_{21}$	X—	—X	—X—
	.		...	
	$S_{32}$	—X	—X—	X—

Figure 9. Overall summary of detection effectiveness data.

Note: some data pertain to only on-line techniques (\*), and some data were collected only in certain phases.

Phase	#Subj.	Measure	Mean	SD	Min.	Max.
1	29	# Faults detected	3.94	1.82	0.00	7.00
1	29	% Faults detected	54.78	26.11	0.00	100.00
1	29(*)	# Faults observable	5.38	1.51	3.00	8.00
1	29(*)	% Faults observable	74.59	20.54	33.33	100.00
1	29(*)	% Detected/observable	70.99	24.01	0.00	100.00
2	13	# Faults detected	3.28	1.96	0.00	7.00
2	13	% Faults detected	39.53	27.25	0.00	100.00
3	32	# Faults detected	4.27	1.86	0.00	8.00
3	32	% Faults detected	49.82	27.44	0.00	100.00
3	32	% Faults felt found	75.10	24.07	0.00	100.00
3	32(*)	# Faults observable	5.61	1.52	3.00	9.00
3	32(*)	% Faults observable	62.11	18.36	25.00	100.00
3	32(*)	% Detected/observable	69.67	27.14	0.00	100.00
3	32(*)	Max. % stmt. covered	97.02	7.83	46.00	100.00
Ave	74	# Faults detected	3.97	1.88	0.00	8.00
Ave	74	% Faults detected	49.96	27.29	0.00	100.00
Ave	61(*)	# Faults observable	5.5	1.5	3.00	9.00
Ave	61(*)	% Faults observable	68.0	20.3	25.0	100.0
Ave	61(*)	% Detected/observable	70.3	25.6	0.0	100.0

Figure 10. Distribution of the number of faults detected broken down by phase. Key: code readers (C), functional testers (F), and structural testers (S).

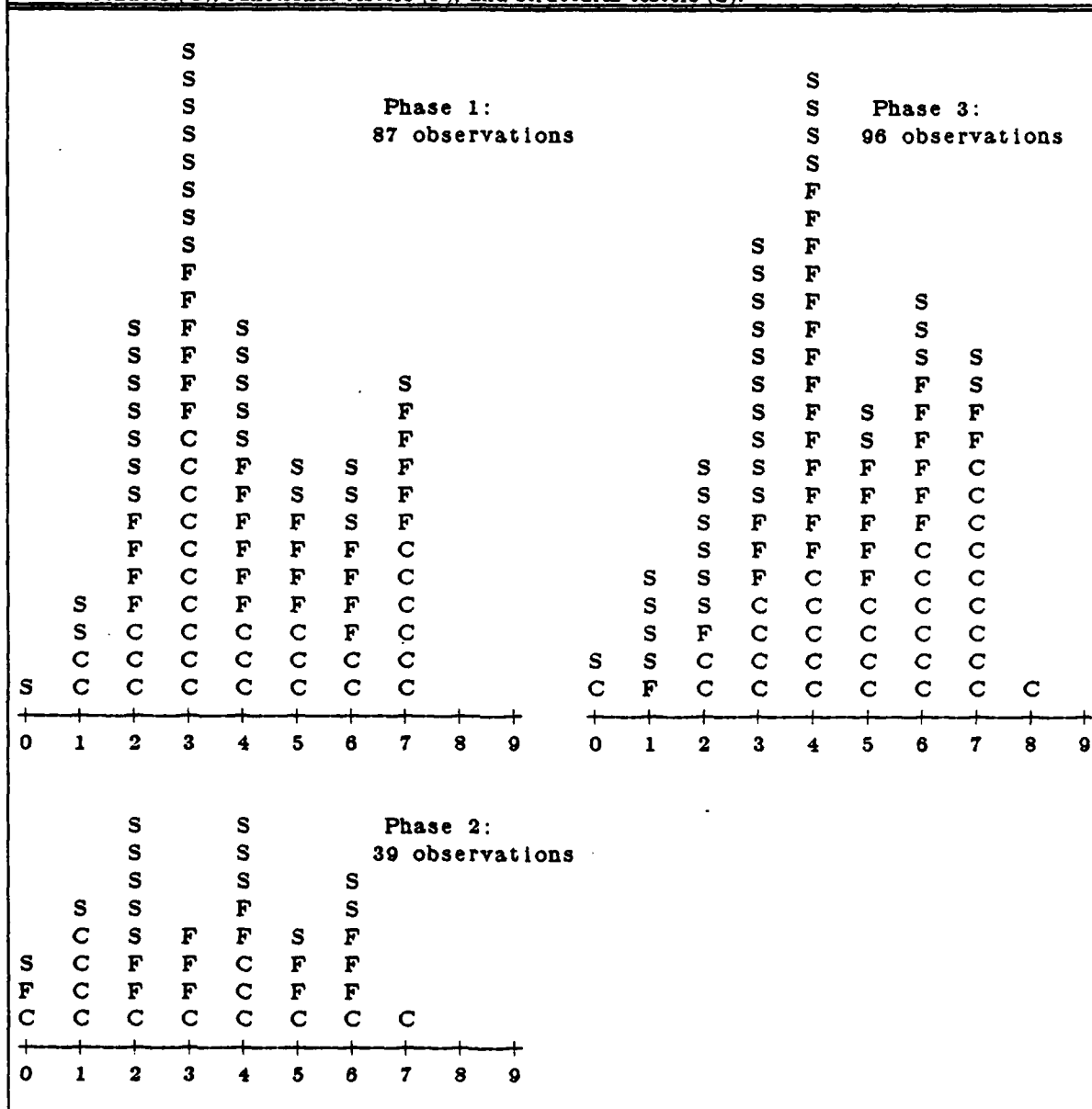
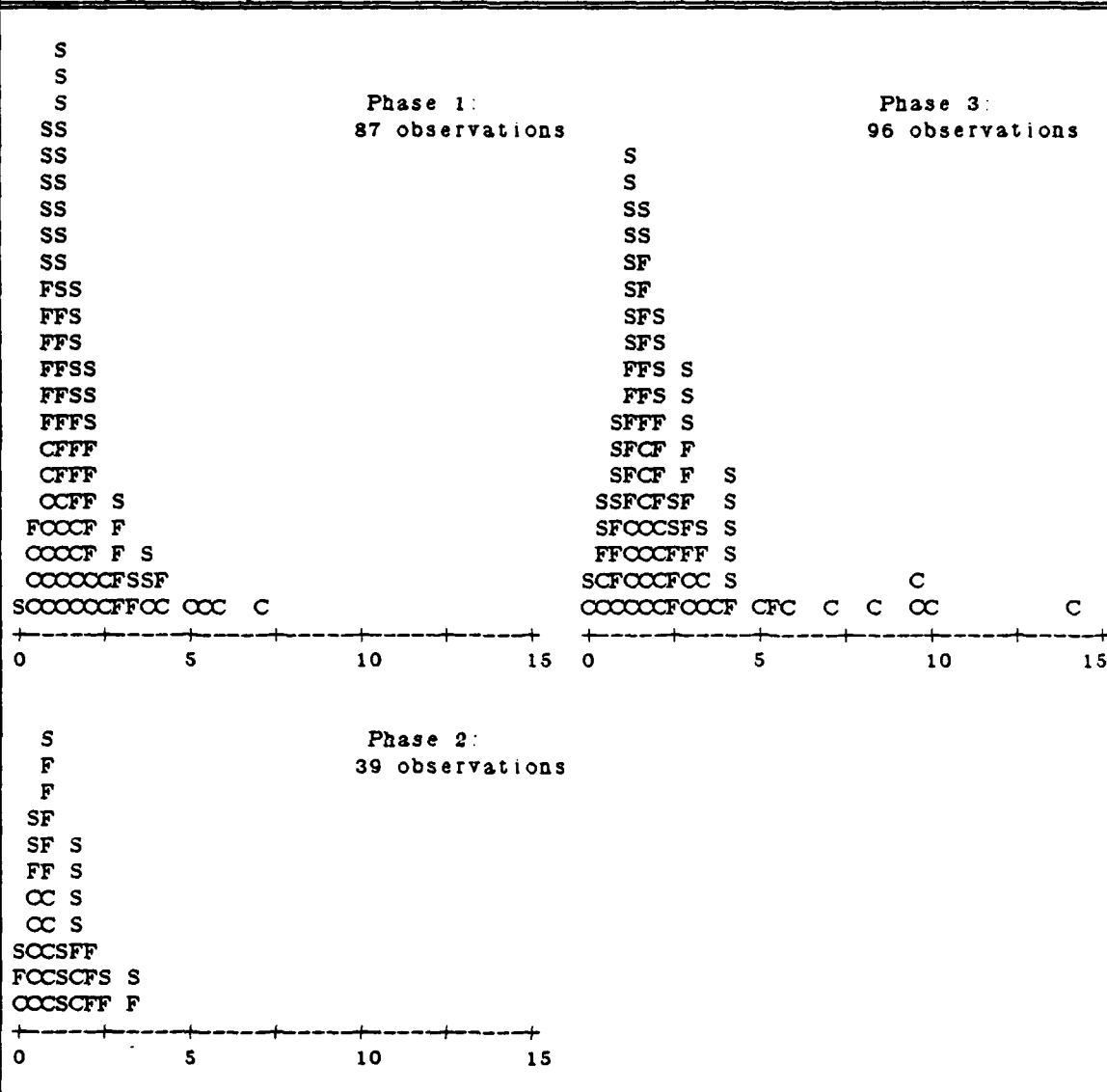


Figure 11. Overall summary for number of faults detected.				
		Phase		
		1	2	3
Effect	Level	Mean(SD)	Mean(SD)	Mean(SD)
Technique	Reading	4.10 (1.93)	3.00 (2.20)	5.09 (1.92)
	Functional	4.45 (1.70)	3.77 (1.83)	4.47 (1.34)
	Structural	3.28 (1.87)	3.08 (1.89)	3.25 (1.80)
Program	Formatter	4.07 (1.62)	3.23 (2.20)	4.19 (1.73)
	Plotter	3.48 (1.45)	3.31 (1.97)	. (.)
	Data type	4.28 (2.25)	. (.)	5.22 (1.75)
	Database	. (.)	3.31 (1.84)	3.41 (1.66)
Expertise	Junior	3.88 (1.89)	3.04 (2.07)	3.90 (1.83)
	Intermed.	4.07 (1.69)	3.83 (1.64)	4.18 (1.99)
	Advanced	. (.)	. (.)	5.00 (1.53)

Figure 12. Overall summary of fault detection cost data. Note: some data pertain to only on-line techniques (*), and some data were collected only in certain phases.						
Phase	#Subj.	Measure	Mean	SD	Mln.	Max.
1	29	# Faults / hour	1.63	1.28	0.00	7.00
1	29	Detection time (hrs)	3.33	2.09	0.75	10.00
2	13	# Faults / hour	0.99	0.81	0.00	3.00
2	13	Detection time (hrs)	4.70	3.02	1.00	14.00
3	32	# Faults / hour	2.33	2.28	0.00	14.00
3	32	Detection time (hrs)	2.75	1.57	0.50	7.25
3	32(*)	Cpu-time (sec)	45.2	56.1	3.0	283.0
3	32(*)	Cpu-time (sec; norm.)	38.5	51.7	2.9	314.4
3	32(*)	Connect time (min)	65.83	50.21	3.50	214.00
3	32(*)	# program runs	5.45	5.00	1.00	24.00
Ave	74	# Faults / hour	1.82	1.80	0.00	14.00
Ave	74	Detection time (hrs)	3.32	2.19	0.50	14.00

Figure 13. Distribution of the fault detection rate (#faults detected per hour) broken down by phase. Key: code readers (C), functional testers (F), and structural testers (S).



**Figure 14. Overall summary for fault detection rate (# faults detected per hour).**

		Phase		
		1	2	3
Effect	Level	Mean(SD)	Mean(SD)	Mean(SD)
Technique	Reading	1.90 (1.83)	0.56 (0.46)	3.33 (3.42)
	Functional	1.58 (0.90)	1.22 (0.91)	1.84 (1.06)
	Structural	1.40 (0.87)	1.18 (0.84)	1.82 (1.24)
Program	Formatter	1.60 (1.39)	0.98 (0.67)	2.15 (1.10)
	Plotter	1.19 (0.83)	0.92 (0.71)	. (.)
	Data type	2.09 (1.42)	. (.)	3.70 (3.26)
	Database	. (.)	1.05 (1.04)	1.14 (0.79)
Expertise	Junior	1.36 (0.97)	1.00 (0.85)	2.14 (2.48)
	Intermed.	2.22 (1.66)	0.96 (0.74)	2.53 (2.48)
	Advanced	. (.)	. (.)	2.36 (1.61)

**Figure 15. Characterization of the faults detected.**

	Code Reading	Functional Testing	Structural Testing	Overall
Omission	55.6 (40.1)	61.0 (39.5)	39.2 (41.6)	52.0 (41.3)
Commission	54.3 (32.1)	53.5 (25.4)	44.3 (26.6)	50.7 (28.4)
Total	54.1 (29.2)	54.6 (24.5)	41.2 (26.1)	50.0 (27.3)
Initial.	64.6 (40.3)	75.0 (36.1)	46.2 (39.8)	61.5 (40.2)
Control	42.8 (36.6)	66.7 (34.9)	48.8 (36.5)	52.8 (37.2)
Data	20.7 (36.6)	28.3 (44.9)	26.8 (41.9)	25.3 (41.0)
Computat.	70.9 (37.0)	64.2 (40.8)	58.8 (43.5)	64.6 (40.6)
Interface	46.7 (38.5)	30.7 (33.5)	24.6 (29.4)	34.1 (35.1)
Cosmetic	16.7 (38.1)	8.3 (28.2)	7.7 (27.2)	10.8 (31.3)
Total	54.1 (29.2)	54.6 (24.5)	41.2 (26.1)	50.0 (27.3)

Figure 18. Characterization of faults detected by the three techniques: 10 omission (O) vs. 24 commission (x). The vertical axis is the number of persons using the particular technique that observed the fault.

	Reading	Functional	Structural
100%	O xxxx Oxxx	OxO Oxxx	Ox x
75%	O x xO xxx xx	x xOxxx x xx	x x O xxx xx x
50%		x Oxx xxxO	xOxxx x Ox x
25%	Ox x Oxx xxxx O	xx xOOxO	xxx xO O OxOO
0%	xOOxO	xOx	Oxx



Figure 17. Characterization of faults detected by the three techniques: Initialization (2-A), computation (8-P), control (7-C), data (3-D), interface (13-I), and cosmetic (1-S). The vertical axis is the number of persons using the particular technique that observed the fault.

	Reading	Functional	Structural
100%	P PIII CIII	PIP PIIC	PC I
75%	P C AP AIP CC	A ACPPC C CC	P C P PCP PC A
50%	CP D	D CPI III	DPIII A CI
25%	DCI SPII I	PI . SIPID	I IIC II D CSDI
0%	CPIID	IDI	PII

Figure 18. Characterization of the faults observable, but not reported.

	Functional Testing	Structural Testing	Overall
Omission	15.7 (25.4)	21.3 (31.8)	18.5 (28.8)
Commission	19.1 (20.0)	20.1 (16.6)	19.6 (18.3)
Total	18.1 (17.8)	19.9 (16.8)	19.0 (17.3)
Initial.	5.0 (15.4)	14.3 (32.2)	9.8 (25.5)
Control	20.3 (30.6)	21.1 (31.4)	20.7 (30.8)
Data	28.6 (43.5)	7.5 (24.5)	18.3 (36.7)
Computat.	16.0 (31.3)	20.1 (37.6)	18.0 (34.5)
Interface	16.1 (20.0)	20.3 (21.5)	18.2 (20.8)
Cosmetic	60.0 (50.3)	85.7 (35.9)	73.2 (44.9)
Total	18.1 (17.8)	19.9 (16.8)	19.0 (17.3)

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER <del>AFOSR</del> TR- 5805	2. GOVT ACCESSION NO. AD-A160136	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) COMPARING THE EFFECTIVENESS OF SOFTWARE TESTING STRATEGIES		5. TYPE OF REPORT & PERIOD COVERED Technical Report
		6. PERFORMING ORG. REPORT NUMBER TR-1501
7. AUTHOR(s) Victor R. Basili and Richard W. Selby, Jr.		8. CONTRACT OR GRANT NUMBER(s) <del>AFOSR</del> -F49620-80-C-001
9. PERFORMING ORGANIZATION NAME AND ADDRESS Department of Computer Science University of Maryland College Park, Maryland 20742		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS G1102F 2304/A2
11. CONTROLLING OFFICE NAME AND ADDRESS Math. & Info. Sciences, AFOSR Bolling AFB Washington, D. C. 20332		12. REPORT DATE May 1985
		13. NUMBER OF PAGES 67
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		15. SECURITY CLASS. (of this report) UNCLASSIFIED
		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report)  Approved for public release; distribution unlimited		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) software testing, functional testing, structural testing, code reading, off-line software review, empirical study, methodology evaluation, software measurement		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) This study compares the strategies of code reading, functional testing, and structural testing in three aspects of software testing: fault detection effectiveness, fault detection cost, and classes of faults detected. Thirty-two professional programmers and 42 advanced students applied the three techniques to four unit-sized programs in a fractional factorial experimental design. The major results of this study are the following. 1) With the professional programmers, code reading detected more software faults and had a higher fault detection		

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

rate than did functional or structural testing, while functional testing detected more faults than did structural testing, but functional and structural testing were not different in fault detection rate. 2) In one advanced student subject group, code reading and functional testing were not different in faults found, but were both superior to structural testing, while in the other advanced student subject group there was no difference among the techniques. 3) With the advanced student subjects, the three techniques were not different in fault detection rate. 4) Number of faults observed, fault detection rate, and total effort in detection depended on the type of software tested. 5) Code reading detected more interface faults than did the other methods. 6) Functional testing detected more control faults than did the other methods. 7) When asked to estimate the percentage of faults detected, code readers gave the most accurate estimates while functional testers gave the least accurate estimates.

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

**END**

**FILMED**

**12-85**

**DTIC**